

# CAN KNOWLEDGE EDITING REALLY CORRECT HALLUCINATIONS?

Anonymous authors

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) suffer from hallucinations, referring to the non-factual information in generated content, despite their superior capacities across tasks. Meanwhile, knowledge editing has been developed as a new popular paradigm to correct erroneous factual knowledge encoded in LLMs with the advantage of avoiding retraining from scratch. However, one common issue of existing evaluation datasets for knowledge editing is that **they do not ensure that LLMs actually generate hallucinated answers to the evaluation questions before editing**. When LLMs are evaluated on such datasets after being edited by different techniques, it is hard to directly adopt the performance to assess the effectiveness of different knowledge editing methods in correcting hallucinations. Thus, the fundamental question remains insufficiently validated: *Can knowledge editing really correct hallucinations in LLMs?* We proposed HalluEditBench to holistically benchmark knowledge editing methods in correcting real-world hallucinations. First, we rigorously construct a massive hallucination dataset with 9 domains, 26 topics and more than 6,000 hallucinations. Then, we evaluate the performance of knowledge editing methods in a holistic way on five dimensions including *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*. Through HalluEditBench, we have provided new insights into the potentials and limitations of different knowledge editing methods in correcting hallucinations, which could inspire future improvements and facilitate progress in the field of knowledge editing. Data and code are available at <https://anonymous.4open.science/r/hallucination-9D6>.

## 1 INTRODUCTION

Large Language Models (LLMs) have shown superior performance in various tasks (Zhao et al., 2023). However, one critical weakness is that they may output hallucinations, referring to the non-factual information in generated content, for reasons such as the limit of models’ internal knowledge scope or fast-changing world facts (Zhang et al., 2023). Considering the high cost of retraining LLMs from scratch, knowledge editing has been designed as a new paradigm to correct erroneous or outdated factual knowledge in LLMs (Wang et al., 2023c).

Although there are many existing question-answering datasets such as WikiData<sub>recent</sub> (Cohen et al., 2024), ZsRE (Yao et al., 2023), and WikiBio (Hartvigsen et al., 2024) widely used for the evaluation of knowledge editing, one common issue is that they do not verify whether LLMs, before applying knowledge editing, actually generate hallucinated answers to the evaluation questions. When such datasets are adopted to evaluate the performance of LLMs after they have been edited, it is hard to directly use the scores to judge the effectiveness of different knowledge editing techniques in correcting hallucinations, which is the motivation to apply knowledge editing to LLMs.

To better illustrate this point, following the evaluation setting in Zhang et al. (2024e), we conducted a preliminary study to examine the pre-edit and post-edit performances of Llama2-7B on the three aforementioned evaluation datasets. As shown in Table 1, we can clearly observe that Llama2-7B achieves a relatively high performance, measured by the rate of answering the evaluation questions

| Method            | WikiData <sub>recent</sub> | ZsRE   | WikiBio |
|-------------------|----------------------------|--------|---------|
| Pre-edit          | 47.40                      | 37.49  | 61.35   |
| Post-edit (ROME)  | 97.37                      | 96.86  | 95.91   |
| Post-edit (MEMIT) | 97.10                      | 95.86  | 94.68   |
| Post-edit (FT-L)  | 56.30                      | 53.82  | 66.70   |
| Post-edit (FT-M)  | 100.00                     | 99.98  | 100.00  |
| Post-edit (LoRA)  | 100.00                     | 100.00 | 100.00  |

Table 1: Performance measured by **Accuracy (%)** of Llama2-7B before editing (“Pre-edit”) and after applying typical knowledge editing methods (“Post-edit”) on common existing evaluation datasets.

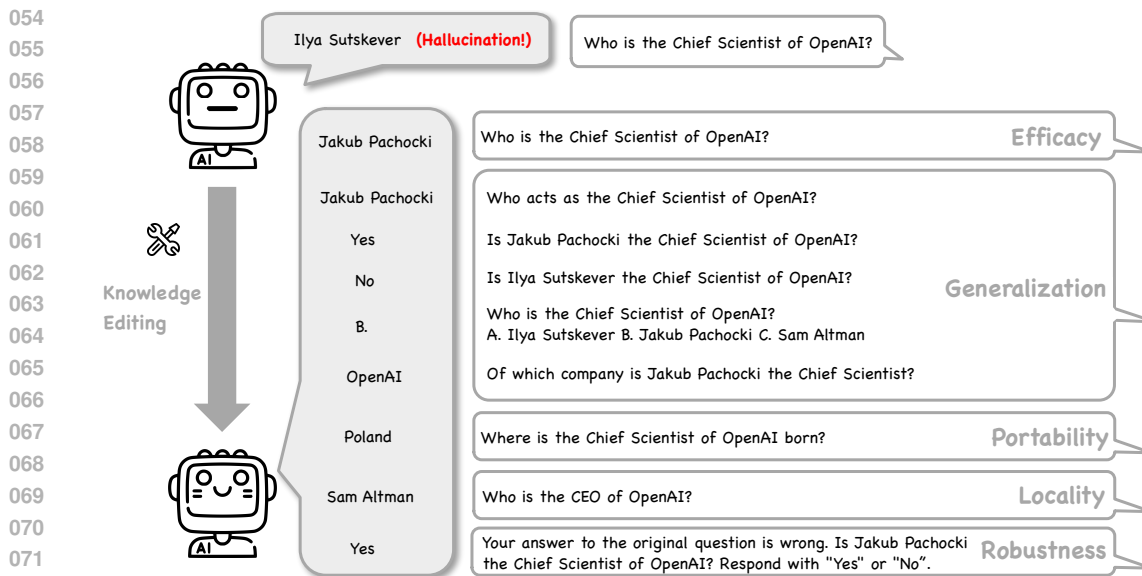


Figure 1: Framework of HalluEditBench. For real-world hallucinations, we holistically assess the performance of knowledge editing on *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*.

correctly (Accuracy (%)), even before applying knowledge editing techniques. Although the knowledge editing methods can bring an increase in accuracy (%), the high post-edit performance on these datasets cannot faithfully reflect the true effectiveness in correcting real-world hallucinations and may cause a distorted assessment. Thus, the fundamental question remains insufficiently validated: **Can knowledge editing really correct hallucinations in LLMs?**

To fill in the essential gap in the field of knowledge editing, we propose HalluEditBench to holistically benchmark knowledge editing techniques to correct real-world hallucinations of LLMs. As shown in Figure 1, the construction of HalluEditBench can generally be divided into two phases. In the first phase, we constructed a massive hallucination dataset encompassing 9 domains and 26 topics based on Wikidata. For each of Llama2-7B, Llama3-8B, and Mistral-v0.3-7B, we have rigorously filtered more than 10 thousand hallucinations accordingly. In the second phase, we sampled around 2,000 hallucinations for each LLM covering all the topics and domains, and then generated evaluation question-answer pairs from five facets including *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*. Through extensive empirical investigation on performance of 7 typical knowledge editing techniques, including FT-L (Zhu et al., 2020; Meng et al., 2022), FT-M (Zhang et al., 2024e), MEMIT (Meng et al., 2023), ROME (Meng et al., 2022), LoRA (Hu et al., 2022), ICE (Zheng et al., 2023), and GRACE (Hartvigsen et al., 2024), regarding the aforementioned five dimensions, we have provided novel insights into their potentials and limitations. A summary of the insights is as follows:

- **The effectiveness of knowledge editing methods in correcting real-world hallucinations could be far from what their performance on existing datasets suggests**, reflecting the potential unreliability of the current assessment of different knowledge editing techniques. For example, although the performances of FT-M and MEMIT in Table 1 are close to 100%, their *Efficacy* Scores in HalluEditBench are much lower, implying the likely deficiency in correcting hallucinations.
- **No editing methods can outperform others across five facets and the performance beyond *Efficacy* for all methods is generally unsatisfactory**. Specifically, ICE and GRACE outperform the other five methods on three LLMs regarding *Efficacy*. All editing methods except ICE only slightly improve or negatively impact the *Generalization* performance. Editing techniques except ICE could even underperform pre-edit LLMs on *Portability*. FT-M and ICE surpass others on *Locality* performance. ICE has a poor *Robustness* performance compared to other methods.
- **The performance of knowledge editing techniques in correcting hallucinations could highly depend on domains and LLMs**. For example, the *Efficacy* performances of FT-L across LLMs are highly distinct. Domains have a large impact on the *Locality* performance of ICE.

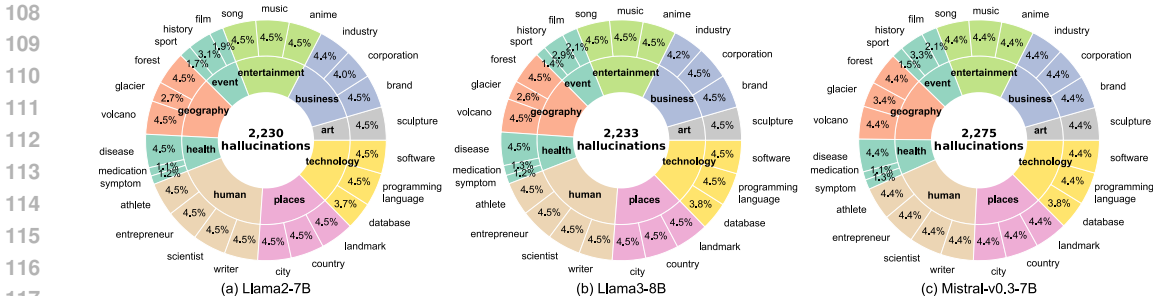


Figure 2: Statistics of HalluEditBench Across Topics and Domains.

## 2 HalluEditBench: HOLISTICALLY BENCHMARKING KNOWLEDGE EDITING METHODS IN CORRECTING REAL-WORLD HALLUCINATIONS

In this section, we will introduce the details of HalluEditBench, including the construction of the massive LLM hallucination dataset, the generation of evaluation question-answering pairs from five dimensions, evaluation metrics and the benchmarked knowledge editing techniques.

### 2.1 HALLUCINATION DATASET CONSTRUCTION

The goal of knowledge editing can generally be defined as transforming existing factual knowledge in the form of a knowledge triplet (subject  $s$ , relation  $r$ , object  $o$ ) into a new one (subject  $s$ , relation  $r$ , object  $o^*$ ). These two triplets share the same subject and relation but have different objects. A knowledge editing operation can be represented as  $e = (s, r, o, o^*)$ . Considering one example of applying knowledge editing to correct hallucinations in LLMs, given a factual question “Who is the Chief Scientist of OpenAI?”, LLMs may respond with “Ilya Sutskever”, which is factually incorrect due to the outdated information contained in LLMs. The editing operation can be  $e = (s = \text{OpenAI}, r = \text{Chief Scientist}, o = \text{Ilya Sutskever}, o^* = \text{Jakub Pachocki})$ . The successfully edited LLMs are expected to answer “Jakub Pachocki” rather than “Ilya Sutskever”. Thus, we need to collect a large scale of knowledge triplets and factual questions to filter hallucinations.

Following existing editing datasets (e.g., WikiData<sub>recent</sub> (Cohen et al., 2024) and WikiBio (Hartvigsen et al., 2024)), we also choose Wikidata as the factual knowledge source. In the *first* step, we retrieved 143, 557 raw knowledge triplets using the Wikidata Query Service (Query date: 2024 September 8) from 26 topics, which can be categorized into 9 domains including *art*, *business*, *entertainment*, *event*, *geography*, *health*, *human*, *places*, and *technology*. Topics were selected based on availability and a minimum of 100 triplets in Wikidata. In the *second* step, we filtered out the triplets that share the same subject and relation while the objects are different, indicating there are more than one answer to questions about the object. When we construct factual questions and compare LLM-generated answers with the triplets, it would be difficult to determine whether LLMs actually hallucinate the questions. For example, for two triplets (Canada, diplomatic relation, India) and (Canada, diplomatic relation, Greece), there are multiple answers to the question “What country has diplomatic relation with Canada?” In the *third* step, following Wang et al. (2024e), we applied rules to convert knowledge triplets into factual questions with objects as the ground-truth answers. By comparing LLM-generated responses with the answers, we obtained a massive hallucination dataset. Specifically, we collected 12, 619, 13, 210, and 14, 366 hallucinations for Llama2-7B, Llama3-8B, and Mistral-v0.3-7B respectively. Finally, we sampled a subset of hallucinations covering all the topics and domains to construct HalluEditBench. The distribution statistics are shown in Figure 2.

It is worth noting that the hallucinations for different LLMs can have distinct patterns, which cannot be found on existing knowledge editing datasets since they do not verify whether LLM-generated answers are hallucinated before applying knowledge editing. **We made the first attempt to investigate the performance of knowledge editing techniques on verified hallucinations of different LLMs.**

### 2.2 EVALUATION QA PAIR GENERATION AND METRICS

After constructing the hallucination dataset, we propose to holistically assess the performance of knowledge editing methods in correcting hallucinations from five facets including *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*. First, we leveraged GPT-4o to generate evaluation question-answering pairs for each facet based on the hallucination dataset as well as the factual

verification questions in Section 2.1. Then we also manually inspect their quality. One example of the evaluation QA pairs for each facet is shown in Figure 1 (More examples are given in Appendix F). The specific prompt design for GPT-4o is shown in Appendix A.

Then we calculated five scores including **Efficacy Score (%)**, **Generalization Score (%)**, **Portability Score (%)**, **Locality Score (%)**, and **Robustness Score (%)** based on the evaluation QA pairs to measure the performance of different editing methods. Except that the Locality Score is defined as the unchanging rate of LLMs’ responses after editing on the Locality evaluation questions, the other scores are calculated by the accuracy of the corresponding evaluation QA pairs. More details are as follows:

**Facet 1: Efficacy** Efficacy Evaluation Questions are the same as factuality verification questions in the hallucination collection to ensure the pre-edit performance is 0 regarding Efficacy Score. Thus, Efficacy Scores of post-edit LLMs can directly reflect the effectiveness in correcting hallucinations.

**Facet 2: Generalization** The Generalization Scores aim to evaluate the ability of LLMs in answering different questions regarding the same knowledge triplet, suggesting the generalization of edited knowledge in diverse scenarios. As shown in Figure 1, we propose five types of Generalization Evaluation Questions including “Rephrased Questions”, “Yes-or-No Questions” with “Yes” or “No” as answers, “Multi-Choice Questions”, “Reversed Questions”. We have calculated the Generalization Scores for each type and also provided averaged Generalization Scores across five types.

**Facet 3: Portability** The Portability Scores aim to measure the ability of LLMs to reason about the downstream effects of edited knowledge. Thus, we design the Efficacy Evaluation Questions with  $N$  hops ( $N = 1 \sim 6$ ) as Portability Evaluation Questions. When  $N = 2$ , the example is shown in Figure 1. When the answer to the question “Who is the Chief Scientist of OpenAI?” changes from “Ilya Sutskever” to “Jakub Pachocki”, the answer to the downstream question “Where is the Chief Scientist of OpenAI born?” should also change from “Russia” to “Poland”.

**Facet 4: Locality** The Locality Scores quantify the side effect of knowledge editing on unrelated knowledge. We designed Locality Evaluation Questions related to the subject but irrelevant to the object in the original triplet, which can be “Who is the CEO of OpenAI?” for the aforementioned example. Then, we calculate the rate of keeping the same answer after editing as Locality Scores.

**Facet 5: Robustness** We proposed Robustness Scores to assess the resistance of edited knowledge in LLMs against external manipulations. Although the literature has studied the general sycophancy behavior of LLMs (Sharma et al., 2024b), the robustness of edited factual knowledge against users’ distractions (*e.g.*, “Your answer to the original question is wrong.”) is underexplored. After post-edit LLMs are tested with Efficacy Evaluation Questions, we further prompted them with Robustness Evaluation Questions, which are exemplified in Figure 1, for  $M$  turns ( $M = 1 \sim 10$ ) and calculated the rate of “Yes” for each round as the Robustness Scores, reflecting the extent to which LLMs insist on corrected knowledge. Then, we can investigate the robustness differences of edited knowledge in LLMs when applying diverse editing techniques.

### 2.3 KNOWLEDGE EDITING TECHNIQUES

We propose to categorize the majority of existing knowledge editing techniques into the following 4 types and chose 7 representative techniques (more details are in Appendix B) in HalluEditBench.

- **Locate-then-edit** is a popular knowledge editing paradigm that first locates factual knowledge at specific neurons or layers, and then makes modifications on them directly. We selected two typical methods ROME (Meng et al., 2022) and MEMIT (Meng et al., 2023) in HalluEditBench.
- **Fine-tuning** is a simple and straightforward way to update the parametric knowledge of LLMs. We selected three variations FT-L (Zhu et al., 2020; Meng et al., 2022), FT-M (Zhang et al., 2024e), and LoRA (Hu et al., 2022), which mitigate the catastrophic forgetting and overfitting issues of standard fine-tuning.
- **In-Context Editing** is a training-free paradigm that associates LLMs with in-context knowledge directly (Zheng et al., 2023; Shi et al., 2024; Fei et al., 2024). We adopted a simple baseline ICE method in Zheng et al. (2023) that puts the new fact in context and does not require demonstrations.
- **Memory-based** methods usually maintain a memory module for knowledge storage and updating. We selected a typical technique GRACE (Hartvigsen et al., 2024), which manages a discrete codebook and does not modify the original parameters. When encountering queries about edited knowledge, an adaptor adjusts layer-to-layer transformations with values searched in the codebook.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

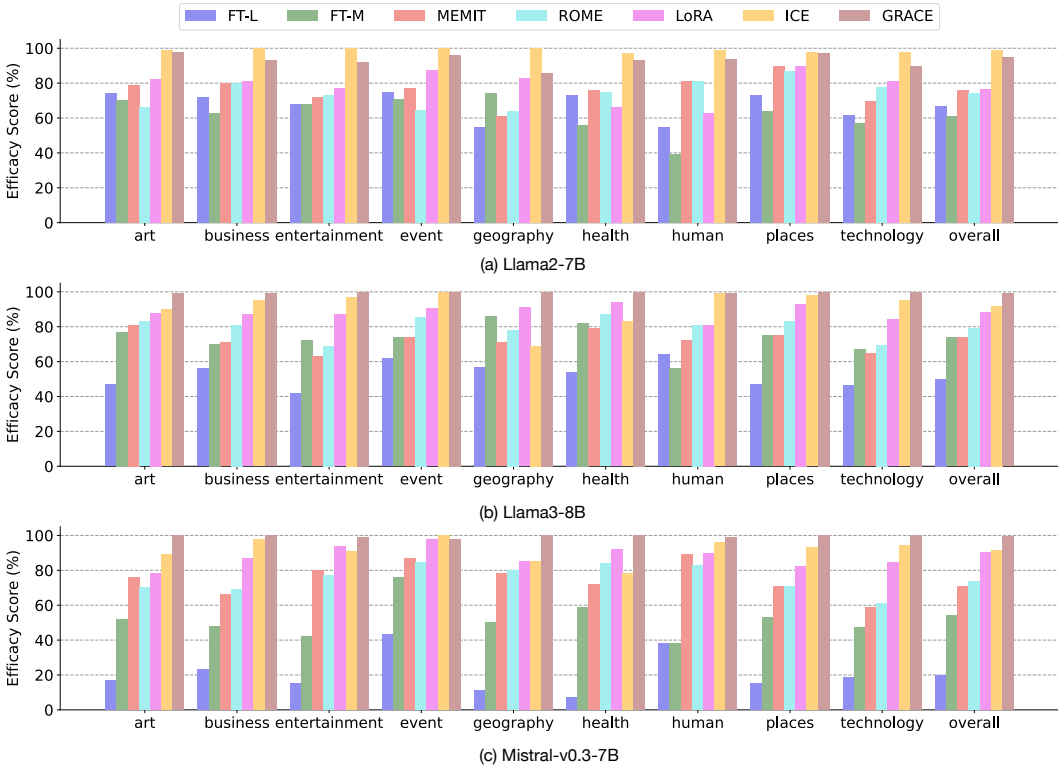


Figure 3: **Efficacy Scores of Knowledge Editing Methods.** The “overall” refers to the Efficacy Score (%) on the whole HalluEditBench embracing 9 domains for different methods. The Efficacy Score on each domain is also reported. Efficacy scores (%) are measured by the accuracy on Efficacy Evaluation Question-answer Pairs, where the pre-edit scores of each LLM are ensured 0.

### 3 RESULTS AND ANALYSIS

In this section, we comprehensively analyze the experiment results on 9 domains and the overall performance on the whole HalluEditBench for different knowledge editing techniques from five facets including *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*.

#### 3.1 FACET 1: EFFICACY

Since we have ensured that LLMs generate hallucinated answers to the Efficacy Evaluation Questions before editing, the pre-edit Efficacy Score for all editing techniques is 0. Thus, Efficacy Scores in Figure 3 can directly reflect the effectiveness of different techniques in correcting real-world hallucinations. We find that **the effectiveness of some techniques can be far from what their performance on previous datasets suggests**, implying the potential unreliability of their previous evaluation. For example, as shown in Table 1, although FT-M achieves near 100% performance in existing datasets such as WikiData<sub>recent</sub>, ZsRE, and WikiBio, its overall Efficacy Scores on Llama2-7B and Mistral-v0.3-7B are only around 60%. There is a similar performance drop for MEMIT.

Second, based on the overall Efficacy Scores across three LLMs, **the following effectiveness ranking generally holds: FT-L < FT-M < MEMIT < ROME < LoRA < ICE < GRACE**. We can observe that ICE and GRACE, which both preserve the original weights in LLMs, outperform the other methods, implying **the potential disadvantage of directly modifying parameters for knowledge editing**.

Third, we notice that **efficacy scores of knowledge editing techniques could highly depend on domains and LLMs**. For example, the scores of FT-L on different domains and LLMs could be highly distinct. The performance of FT-L and FT-M on Llama3-8B is higher than that on Mistral-v0.3-7B.

**Insight 1:** (1) The current assessment of knowledge editing could be unreliable; (2) ICE and GRACE outperform parameter-modifying editing techniques such as fine-tuning and “Locate-then-Edit” methods on *Efficacy*; (3) Domains and LLMs could have a high impact on *Efficacy*.

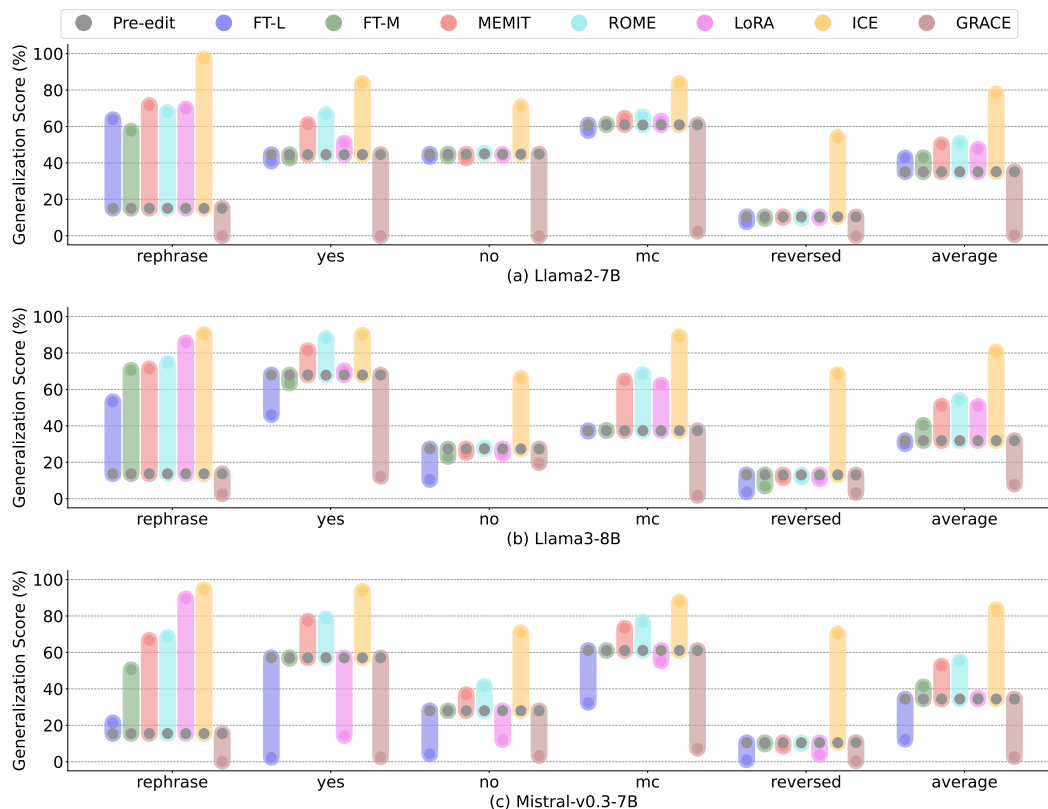


Figure 4: **Generalization Scores of Knowledge Editing Methods.** Generalization Scores (%) are measured by accuracy on five types of Generalization Evaluation Questions including Rephrased Questions (“rephrase”), Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), Reversed Questions (“reversed”). The “average” refers to averaged scores over five question types. The figure only shows the overall Generalization Scores for each type on the whole HalluEditBench. Generalization Scores for each domain are given in Appendix E.1.

### 3.2 FACET 2: GENERALIZATION

As shown in Figure 4, even though the pre-edit Efficacy Score performances for different editing techniques on three LLMs are ensured 0, it is worth noting that the pre-edit Generalization Score performance is not 0 regarding each question type, illustrating that **the manifestation of hallucination actually depends on the design of question prompts**. Given a group of diverse question prompts for the same knowledge triplet, LLMs may hallucinate some questions but answer others correctly.

Surprisingly, we find that **post-edit Generalization Scores could even be lower than pre-edit scores** for the same LLM and question type, demonstrating the potential negative effect caused by knowledge editing. In more detail, we can observe a clear performance drop for GRACE across all the question types, and for FT-L and LoRA on some question types.

Comparing the ranking of Efficacy Scores in Figure 3 with Figure 4, we can explicitly see that **higher Efficacy Scores do not also necessarily indicate higher Generalization Scores**. Especially, although GRACE almost surpasses all the other editing techniques regarding Efficacy Scores, it largely degrades the Generalization Scores compared to pre-edit performance. In addition, **all editing methods except ICE only slightly improve or even hurt Generalization Scores**.

**Insight 2:** (1) The manifestation of hallucination depends on question design; (2) Higher *Efficacy* Scores do not also necessarily indicate higher *Generalization* Scores; (3) All editing techniques except ICE only slightly improve or negatively impact the *Generalization* performance.

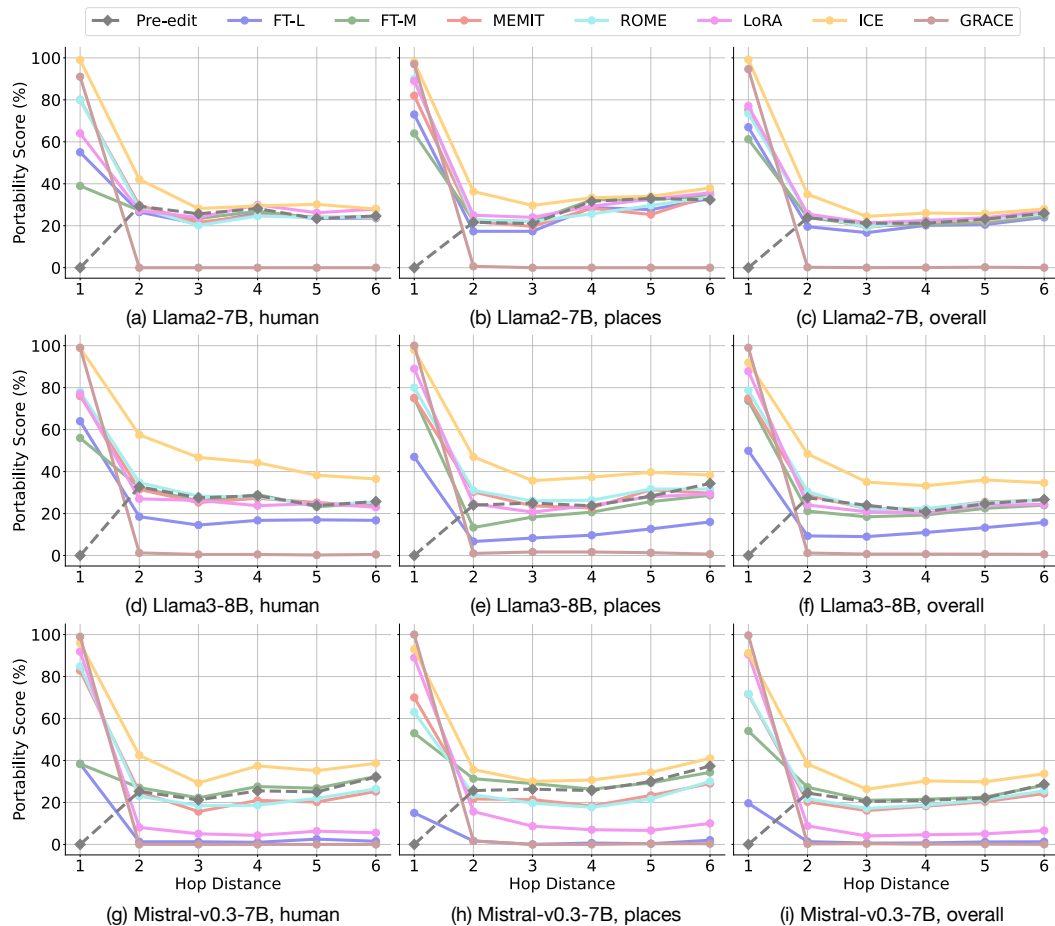


Figure 5: **Portability Scores of Knowledge Editing Methods.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions with  $N$  hops ( $N = 1 \sim 6$ ). The Portability Evaluation Questions are the same as Efficacy Evaluation Questions when  $N$  is 1. The Portability Scores on two domains “human” and “places” are reported in the figure. The results for more domains are given in Appendix E.2. The “overall” refers to the Portability Score (%) on the whole HalluEditBench embracing 9 domains.

### 3.3 FACET 3: PORTABILITY

Figure 5 demonstrates the pre-edit and post-edit Portability Scores for Portability Evaluation Questions with  $N$  hops ( $N = 1 \sim 6$ ). When  $N = 1$ , the Portability Evaluation Questions are the same as Efficacy Evaluation Questions, suggesting that the Portability Scores are 0. Similar to Figure 4, we discover that the pre-edit Portability Scores are not zero for  $2 \sim 6$  hops, indicating **LLMs do not necessarily need to reason based on single-hop knowledge to answer multi-hop questions**. We hypothesize that this is because LLMs may directly memorize the answers to multi-hop questions.

We surprisingly find that except that ICE may bring marginal improvement to the pre-edit performance, **the other knowledge editing techniques even mostly underperform pre-edit Portability Scores**, showing another type of negative effect of knowledge editing and **LLMs may not really reason with the edited knowledge in multi-hop questions** for most knowledge editing methods. Comparing single-hop and multi-hop performance, we observe a sharp decrease for all the editing methods, which further underscores **the challenges of answering multi-hop questions with edited knowledge**.

**Insight 3:** (1) LLMs may memorize answers rather than reason based on single-hop knowledge for multi-hop questions; (2) Editing methods marginally improve or degrade pre-edit Portability Scores, implying LLMs may not really reason with edited knowledge in multi-hop questions.

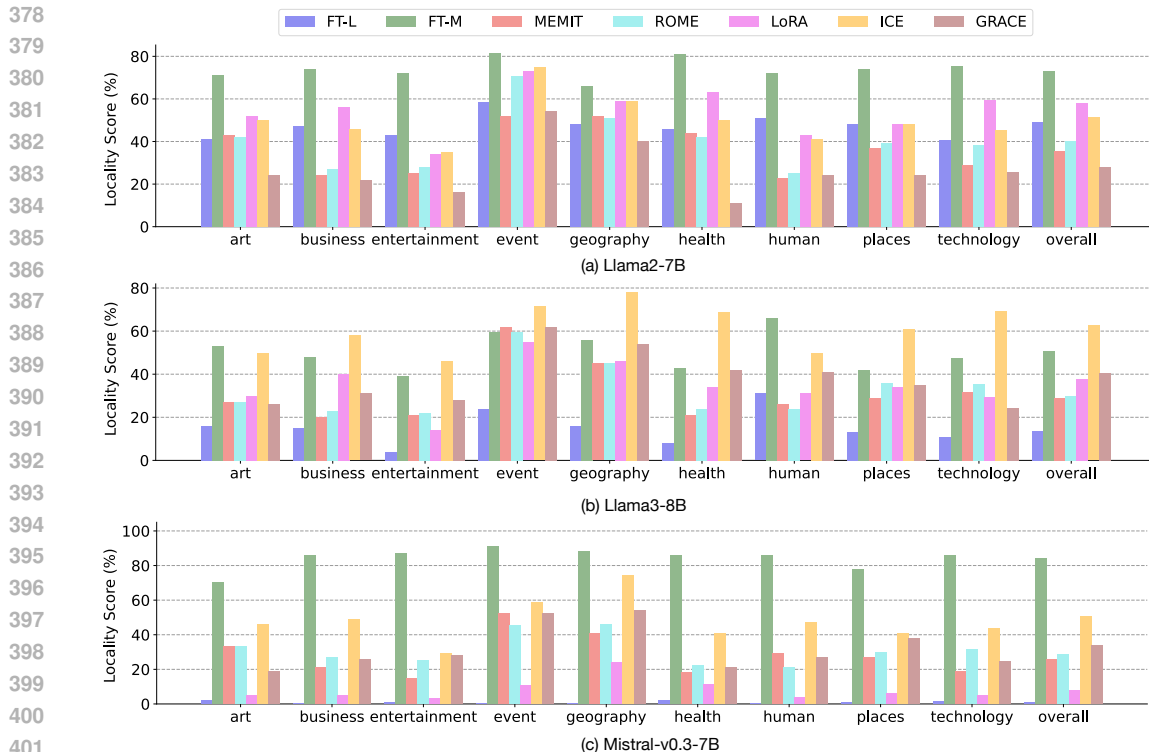


Figure 6: **Locality Scores of Knowledge Editing Methods.** Locality Scores (%) are measured by the unchanging rate on Locality Evaluation Questions after applying knowledge editing methods on LLMs. A higher Locality Score indicates that there is a higher percentage of LLMs’ answers to the unrelated questions keeping the same and a less side effect on general knowledge in LLMs. The “overall” refers to the Locality Score (%) on the whole HalluEditBench embracing 9 domains for different methods. The Locality Score on each domain is also reported in the figure.

### 3.4 FACET 4: LOCALITY

Figure 6 shows the Locality Scores of different editing techniques in each domain and the whole HalluEditBench, reflecting the side effect of knowledge editing on unrelated knowledge encoded in LLMs. Based on the overall Locality Scores, we can observe that **the performance of all editing methods except FT-M and ICE is unsatisfactory**. In particular, the overall Locality Scores for all editing techniques except FT-M and ICE on Llama3-8B and Mistral-v0.3-7B are below 40%, suggesting a high undesired impact on LLMs’ answers to unrelated factual questions, though FT-M achieves an overall score of around 80% on Mistral-v0.3-7B and ICE gains 60% on Llama3-8B.

Furthermore, we notice that **domains and LLMs have a high impact on the Locality Scores of knowledge editing methods**. For example, the Locality Score for ICE in the geography domain in Llama3-8B is around 80%, while the performance drops to only about 40% in the entertainment domain for the same LLM. Although FT-M obtains a Locality Score of around 80% in the entertainment domain on Mistral-v0.3-7B, its performance in the same domain on Llama3-8B is below 40%.

Due to the impact of LLMs, we observe that **the rankings by Locality Scores for editing techniques on different LLMs are highly distinct**. For example, the Locality ranking on Llama2-7B is GRACE < MEMIT < ROME < FT-L < ICE < LoRA < FT-M. However, the ranking changes to FT-L < LoRA < MEMIT < ROME < GRACE < ICE < FT-M on Mistral-v0.3-7B. Comparing Figure 3 with Figure 6, we find **there is no noticeable correlation between Efficacy and Locality for different editing techniques**. FT-M achieves relatively high Locality Scores despite its low Efficacy Scores.

**Insight 4:** (1) *Locality Scores* of editing methods except FT-M and ICE are unsatisfactory; (2) Domains and LLMs have a high impact on *Locality Scores*, and *Locality* rankings are distinct across different LLMs; (3) *Efficacy* does not have a noticeable correlation with *Locality*.



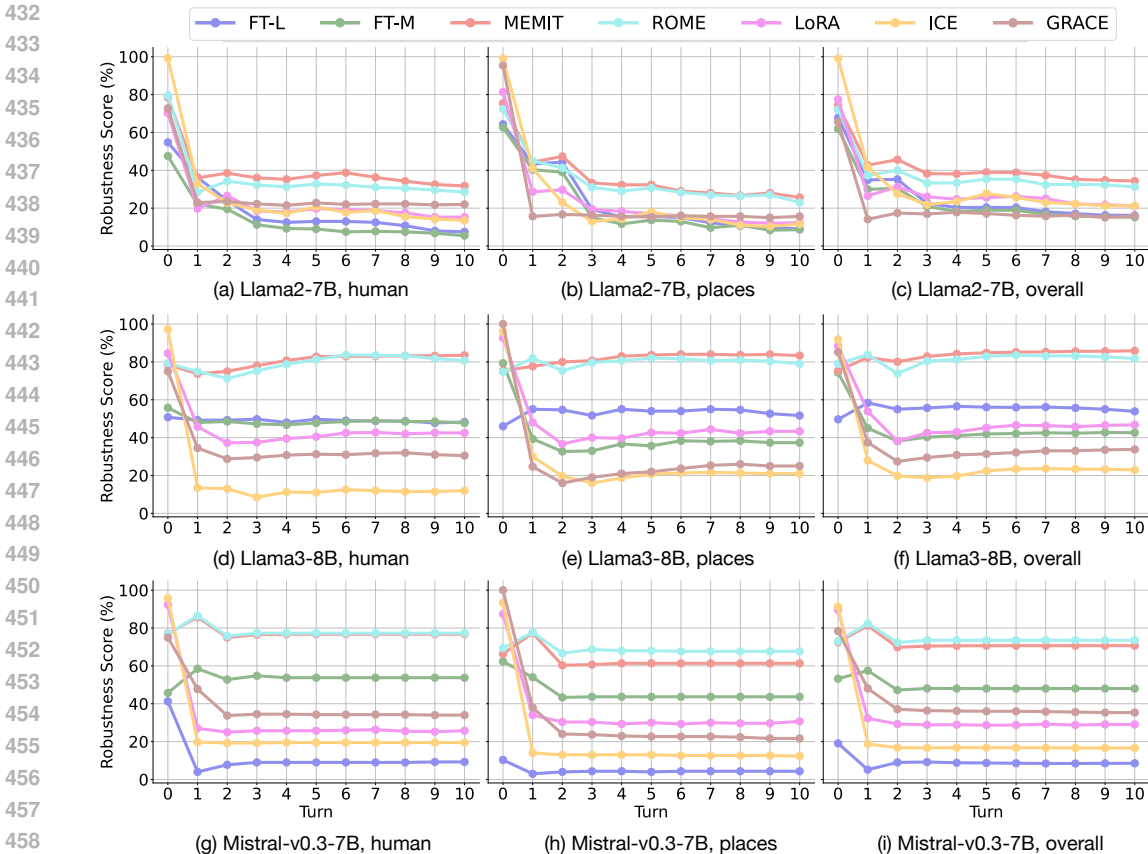


Figure 7: **Robustness Scores of Knowledge Editing Methods.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with  $M$  turns ( $M = 1 \sim 10$ ). We regard Efficacy Scores as the Robustness Scores when  $M$  is 0. The Robustness Scores on two domains “human” and “places” are reported in the figure. The results for more domains are given in Appendix E.3. The “overall” refers to the Robustness Score (%) on the whole HalluEditBench embracing 9 domains.

### 3.5 FACET 5: ROBUSTNESS

We proposed Robustness Scores (%) to evaluate the resistance of edited knowledge against distractions in prompts. Initially ( $M = 0$ ), LLMs are assessed with Efficacy Evaluation Questions. Then ( $M = 1 \sim 10$ ), LLMs are sequentially prompted with Robustness Evaluation Questions, which are exemplified in Figure 1, for  $M$  turns. Robustness Scores are calculated with the percentage of “Yes” in each round. A higher Robustness Score indicates that there is a larger percentage of LLMs can resist external manipulations in the prompt and a higher extent of robustness for the edited knowledge.

First, based on overall Robustness Scores, we observe that **LLMs themselves have a large impact on the robustness of edited knowledge. The same editing method could show distinct trends as turns increase on different LLMs.** For example, all editing methods have a sharp drop when turns go up on Llama2-7B, showing a low level of robustness. However, MEMIT, ROME on Llama3-8B and Mistral-v0.3-7B maintain almost the same and relatively high performance as turns increase, suggesting a comparatively high level of robustness for the edited knowledge.

Then, we notice that **both ICE and GRACE have a low level of robustness** though they outperform the other five editing techniques regarding Efficacy Scores, showing **the potential weaknesses on robustness of parameter-preserving knowledge editing methods.** However, parameter-modifying editing techniques do not necessarily have high robustness, which is exemplified by LoRA.

**Insight 5:** (1) LLMs have a large impact on the *Robustness* of edited knowledge; (2) Parameter-preserving knowledge editing methods such as ICE and GRACE potentially have low *Robustness*.

## 4 RELATED WORK

Knowledge editing techniques have attracted increasing attention for their efficiency advantages in addressing obsolete or hallucinated information in LLMs (Wang et al., 2023c; Zhang et al., 2024e). In general, the existing editing techniques can be categorized into four types including *Locate-then-edit* (Meng et al., 2022; 2023), *Fine-tuning based* (Gangadhar & Stratos, 2024; Zhu et al., 2020; Wang et al., 2024a), *In-Context Editing* (Zheng et al., 2023; Shi et al., 2024; Fei et al., 2024), and *Memory-based* (Wang et al., 2024d; Hartvigsen et al., 2024; Mitchell et al., 2022; Yu et al., 2023). Recently, many benchmarks have been built to investigate the properties of knowledge editing from different perspectives (Rosati et al., 2024; Wu et al., 2023; Ge et al., 2024a; Ma et al., 2023; Wei et al., 2023; 2024a; Zhong et al., 2023; Lin et al., 2024; Huang et al., 2024c; Liu et al., 2024c; Akyürek et al., 2023; Li et al., 2024a;d; 2023b; Gu et al., 2024; Powell et al., 2024). For example, Gu et al. (2024) proposed a benchmark to assess the side effect of 4 popular editing methods on 3 LLMs across 8 general capacity tasks. Rosati et al. (2024) built a new evaluation protocol to measure the efficacy and impact of knowledge editing in long-form generation. Wei et al. (2024a) introduced a multilingual knowledge editing benchmark embracing five languages. However, considering the fundamental motivation of applying knowledge editing to LLMs, which is to correct hallucinations, there is a pressing need to build a real-world hallucination dataset with rigorous verification and systematically analyze the performance of different editing methods. Thus, we proposed HalluEditBench to fill in the gap and provided new insights to facilitate the progress in the field of knowledge editing.

## 5 CONCLUSION

In this paper, we have built a new benchmark HalluEditBench to holistically assess diverse knowledge editing techniques in correcting real-world hallucinations. First, we meticulously construct a comprehensive hallucination dataset based on Wikidata with 9 domains, 26 topics, and more than 6,000 hallucinations. Then, we systematically investigate the performance of different knowledge editing methods from five perspectives including *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*. Our findings reveal a disconnect between current benchmarks and real-world performance. For example, methods like FT-M and MEMIT, while achieving near-perfect scores on existing datasets, perform poorly in HalluEditBench. No method consistently outperforms others across all facets: ICE and GRACE excel in Efficacy, but ICE lags in Robustness, and most methods show limited improvement or even degradation in Generalization and Portability. These results indicate that editing methods often fail to enable true reasoning. This study highlights the limitations of current techniques and benchmarks, emphasizing the need for more reliable approaches to address hallucinations effectively. Our findings offer actionable insights to inspire future advancements in knowledge editing for large language models.

## REFERENCES

- 540  
541  
542 Zhila Aghajari, Eric PS Baumer, and Dominic DiFranzo. Reviewing interventions to address  
543 misinformation: the need to expand our vision beyond an individualistic focus. *Proceedings of the*  
544 *ACM on Human-Computer Interaction*, 7(CSCW1):1–34, 2023.
- 545 Afra Feyza Akyürek, Eric Pan, Garry Kuwanto, and Derry Wijaya. Dune: Dataset for unified editing.  
546 *ArXiv preprint*, abs/2311.16087, 2023. URL <https://arxiv.org/abs/2311.16087>.
- 547  
548 Joseph B Bak-Coleman, Ian Kennedy, Morgan Wack, Andrew Beers, Joseph S Schafer, Emma S  
549 Spiro, Kate Starbird, and Jevin D West. Combining interventions to reduce the spread of viral  
550 misinformation. *Nature Human Behaviour*, 6(10):1372–1380, 2022.
- 551 Alimohammad Beigi, Zhen Tan, Nivedh Mudiam, Canyu Chen, Kai Shu, and Huan Liu. Model  
552 attribution in machine-generated disinformation: A domain generalization approach with super-  
553 vised contrastive learning. *ArXiv preprint*, abs/2407.21264, 2024. URL <https://arxiv.org/abs/2407.21264>.
- 554  
555 Baolong Bi, Shenghua Liu, Lingrui Mei, Yiwei Wang, Pengliang Ji, and Xueqi Cheng. Decod-  
556 ing by contrasting knowledge: Enhancing llms’ confidence on edited facts. *ArXiv preprint*,  
557 abs/2405.11613, 2024a. URL <https://arxiv.org/abs/2405.11613>.
- 558  
559 Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, Hongcheng Gao, Junfeng Fang, and Xueqi  
560 Cheng. Struedit: Structured outputs enable the fast and accurate knowledge editing for large  
561 language models. *ArXiv preprint*, abs/2409.10132, 2024b. URL <https://arxiv.org/abs/2409.10132>.
- 562  
563 Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, Hongcheng Gao, Yilong Xu, and Xueqi  
564 Cheng. Adaptive token biaser: Knowledge editing via biasing key entities. *arXiv preprint arXiv:*  
565 *2406.12468*, 2024c.
- 566  
567 Yuchen Cai, Ding Cao, Rongxi Guo, Yaqin Wen, Guiquan Liu, and Enhong Chen. Editing knowledge  
568 representation of language model via rephrased prefix prompts. *ArXiv preprint*, abs/2403.14381,  
569 2024a. URL <https://arxiv.org/abs/2403.14381>.
- 570  
571 Yuchen Cai, Ding Cao, Rongxi Guo, Yaqin Wen, Guiquan Liu, and Enhong Chen. Locating and  
572 mitigating gender bias in large language models. *ArXiv preprint*, abs/2403.14409, 2024b. URL  
573 <https://arxiv.org/abs/2403.14409>.
- 574  
575 Canyu Chen and Kai Shu. Combating misinformation in the age of llms: Opportunities and challenges.  
576 *AI Magazine*, 2024a. doi: 10.1002/aaai.12188. URL <https://doi.org/10.1002/aaai.12188>.
- 577  
578 Canyu Chen and Kai Shu. Can LLM-generated misinformation be detected? In *The Twelfth*  
579 *International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=ccxD4mtkTU>.
- 580  
581 Canyu Chen, Haoran Wang, Matthew Shapiro, Yunyu Xiao, Fei Wang, and Kai Shu. Combating  
582 health misinformation in social media: Characterization, detection, intervention, and open issues.  
583 *ArXiv preprint*, abs/2211.05289, 2022. URL <https://arxiv.org/abs/2211.05289>.
- 584  
585 Canyu Chen, Baixiang Huang, Zekun Li, Zhaorun Chen, Shiyang Lai, Xiong Xiao Xu, Jia-Chen Gu,  
586 Jindong Gu, Huaxiu Yao, Chaowei Xiao, Xifeng Yan, William Yang Wang, Philip Torr, Dawn  
587 Song, and Kai Shu. Can editing llms inject harm? *ArXiv preprint*, abs/2407.20224, 2024a. URL  
588 <https://arxiv.org/abs/2407.20224>.
- 589  
590 Qizhou Chen, Taolin Zhang, Dongyang Li, Longtao Huang, Hui Xue, Chengyu Wang, and Xiaofeng  
591 He. Lifelong knowledge editing for llms with retrieval-augmented continuous prompt learning.  
592 *ArXiv preprint*, abs/2405.03279, 2024b. URL <https://arxiv.org/abs/2405.03279>.
- 593  
Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Journey to the center of the  
knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate  
knowledge neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38,  
pp. 17817–17825, 2024c.

- 594 Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Knowledge localization: Mission  
595 not accomplished? enter query localization! *ArXiv preprint*, abs/2405.14117, 2024d. URL  
596 <https://arxiv.org/abs/2405.14117>.  
597
- 598 Keyuan Cheng, Muhammad Asif Ali, Shu Yang, Gang Ling, Yuxuan Zhai, Haoyang Fei, Ke Xu,  
599 Lu Yu, Lijie Hu, and Di Wang. Leveraging logical rules in knowledge editing: A cherry on the top.  
600 *ArXiv preprint*, abs/2405.15452, 2024a. URL <https://arxiv.org/abs/2405.15452>.  
601
- 602 Keyuan Cheng, Gang Lin, Haoyang Fei, Lu Yu, Muhammad Asif Ali, Lijie Hu, Di Wang, et al.  
603 Multi-hop question answering under temporal knowledge editing. *ArXiv preprint*, abs/2404.00492,  
604 2024b. URL <https://arxiv.org/abs/2404.00492>.  
605
- 606 Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects  
607 of knowledge editing in language models. *Transactions of the Association for Computational  
608 Linguistics*, 12:283–298, 2024.
- 609 Jingcheng Deng, Zihao Wei, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. Unke:  
610 Unstructured knowledge editing in large language models. *ArXiv preprint*, abs/2405.15349, 2024.  
611 URL <https://arxiv.org/abs/2405.15349>.
- 612 Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Xiang Wang, Xiangnan He, and Tat-seng  
613 Chua. Alphaedit: Null-space constrained knowledge editing for language models. *ArXiv preprint*,  
614 abs/2410.02355, 2024. URL <https://arxiv.org/abs/2410.02355>.  
615
- 616 Weizhi Fei, Xueyan Niu, Guoqing Xie, Yanhua Zhang, Bo Bai, Lei Deng, and Wei Han. Re-  
617 trieval meets reasoning: Dynamic in-context editing for long-text understanding. *ArXiv preprint*,  
618 abs/2406.12331, 2024. URL <https://arxiv.org/abs/2406.12331>.
- 619 Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. A primer on the inner  
620 workings of transformer-based language models. *ArXiv preprint*, abs/2405.00208, 2024. URL  
621 <https://arxiv.org/abs/2405.00208>.  
622
- 623 Govind Gangadhar and Karl Stratos. Model editing by pure fine-tuning. *ArXiv preprint*,  
624 abs/2402.11078, 2024. URL <https://arxiv.org/abs/2402.11078>.  
625
- 626 Huaizhi Ge, Frank Rudzicz, and Zining Zhu. How well can knowledge edit methods edit perplexing  
627 knowledge? *ArXiv preprint*, abs/2406.17253, 2024a. URL [https://arxiv.org/abs/2406.  
628 17253](https://arxiv.org/abs/2406.17253).
- 629 Xiou Ge, Ali Mousavi, Edouard Grave, Armand Joulin, Kun Qian, Benjamin Han, Mostafa Arefiyan,  
630 and Yunyao Li. Time sensitive knowledge editing through efficient finetuning. *ArXiv preprint*,  
631 abs/2406.04496, 2024b. URL <https://arxiv.org/abs/2406.04496>.  
632
- 633 Hengrui Gu, Kaixiong Zhou, Xiaotian Han, Ninghao Liu, Ruobing Wang, and Xin Wang.  
634 Pokemqa: Programmable knowledge editing for multi-hop question answering. *ArXiv preprint*,  
635 abs/2312.15194, 2023. URL <https://arxiv.org/abs/2312.15194>.
- 636 Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun  
637 Peng. Model editing harms general abilities of large language models: Regularization to the rescue.  
638 *ArXiv preprint*, abs/2401.04700, 2024. URL <https://arxiv.org/abs/2401.04700>.  
639
- 640 Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. Model editing at scale leads to gradual and  
641 catastrophic forgetting. *ArXiv preprint*, abs/2401.07453, 2024. URL [https://arxiv.org/abs/  
2401.07453](https://arxiv.org/abs/<br/>642 2401.07453).
- 643 Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi.  
644 Aging with grace: Lifelong model editing with discrete key-value adaptors. *Advances in Neural  
645 Information Processing Systems*, 36, 2024.  
646
- 647 Katrin Hartwig, Frederic Doell, and Christian Reuter. The landscape of user-centered misinformation  
interventions-a systematic literature review. *ACM Computing Surveys*, 56(11):1–36, 2024.

- 648 Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. Does localization inform editing?  
649 surprising differences in causality-based localization vs. knowledge editing in language models.  
650 *Advances in Neural Information Processing Systems*, 36, 2024a.
- 651 Peter Hase, Thomas Hofweber, Xiang Zhou, Elias Stengel-Eskin, and Mohit Bansal. Fundamental  
652 problems with model editing: How should rational belief revision work in llms? *ArXiv preprint*,  
653 abs/2406.19354, 2024b. URL <https://arxiv.org/abs/2406.19354>.
- 654 Bing He, Mustaque Ahamad, and Srijan Kumar. Reinforcement learning-based counter-  
655 misinformation response generation: a case study of covid-19 vaccine misinformation. In *Proceed-*  
656 *ings of the ACM Web Conference 2023*, pp. 2698–2709, 2023.
- 657 Jason Hoelscher-Obermaier, Julia Persson, Esben Kran, Ioannis Konstas, and Fazl Barez. Detecting  
658 edit failures in large language models: An improved specificity benchmark. *ArXiv preprint*,  
659 abs/2305.17553, 2023. URL <https://arxiv.org/abs/2305.17553>.
- 660 Cheng-Hsun Hsueh, Paul Kuo-Ming Huang, Tzu-Han Lin, Che-Wei Liao, Hung-Chieh Fang, Chao-  
661 Wei Huang, and Yun-Nung Chen. Editing the mind of giants: An in-depth exploration of pitfalls  
662 of knowledge editing in large language models. *ArXiv preprint*, abs/2406.01436, 2024. URL  
663 <https://arxiv.org/abs/2406.01436>.
- 664 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
665 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth Interna-*  
666 *tional Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.  
667 OpenReview.net, 2022. URL <https://openreview.net/forum?id=nZevKeeFYf9>.
- 668 Wenyue Hua, Jiang Guo, Mingwen Dong, Henghui Zhu, Patrick Ng, and Zhiguo Wang. Propagation  
669 and pitfalls: Reasoning-based assessment of knowledge editing through counterfactual tasks. *ArXiv*  
670 *preprint*, abs/2401.17585, 2024. URL <https://arxiv.org/abs/2401.17585>.
- 671 Baixiang Huang, Canyu Chen, and Kai Shu. Authorship attribution in the era of llms: Problems,  
672 methodologies, and challenges. *ArXiv preprint*, abs/2408.08946, 2024a. URL <https://arxiv.org/abs/2408.08946>.
- 673 Baixiang Huang, Canyu Chen, and Kai Shu. Can large language models identify authorship?, 2024b.  
674 URL <https://arxiv.org/abs/2403.08213>.
- 675 Han Huang, Haitian Zhong, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. Kebench: A benchmark  
676 on knowledge editing for large vision-language models. *ArXiv preprint*, abs/2403.07350, 2024c.  
677 URL <https://arxiv.org/abs/2403.07350>.
- 678 Houcheng Jiang, Junfeng Fang, Tianyu Zhang, An Zhang, Ruipeng Wang, Tao Liang, and Xiang  
679 Wang. Neuron-level sequential editing for large language models. *ArXiv preprint*, abs/2410.04045,  
680 2024a. URL <https://arxiv.org/abs/2410.04045>.
- 681 Yuxin Jiang, Yufei Wang, Chuhan Wu, Wanjun Zhong, Xingshan Zeng, Jiahui Gao, Liangyou Li, Xin  
682 Jiang, Lifeng Shang, Ruiming Tang, et al. Learning to edit: Aligning llms with knowledge editing.  
683 *ArXiv preprint*, abs/2402.11905, 2024b. URL <https://arxiv.org/abs/2402.11905>.
- 684 Jiaqi Li, Miaozen Du, Chuanyi Zhang, Yongrui Chen, Nan Hu, Guilin Qi, Haiyun Jiang, Siyuan  
685 Cheng, and Bozhong Tian. Mike: A new benchmark for fine-grained multimodal entity knowledge  
686 editing. *ArXiv preprint*, abs/2402.14835, 2024a. URL <https://arxiv.org/abs/2402.14835>.
- 687 Shuaiyi Li, Yang Deng, Deng Cai, Hongyuan Lu, Liang Chen, and Wai Lam. Consecutive model  
688 editing with batch alongside hook layers. *ArXiv preprint*, abs/2403.05330, 2024b. URL <https://arxiv.org/abs/2403.05330>.
- 689 Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. Pmet: Precise model editing  
690 in a transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp.  
691 18564–18572, 2024c.
- 692 Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. Unveiling the  
693 pitfalls of knowledge editing for large language models. *ArXiv preprint*, abs/2310.02129, 2023a.  
694 URL <https://arxiv.org/abs/2310.02129>.

- 702 Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. Unveiling the  
703 pitfalls of knowledge editing for large language models. In *The Twelfth International Conference*  
704 *on Learning Representations*, 2024d. URL <https://openreview.net/forum?id=fNktD3ib16>.  
705
- 706 Zichao Li, Ines Arous, Siva Reddy, and Jackie Chi Kit Cheung. Evaluating dependencies in fact  
707 editing for language models: Specificity and implication awareness. In *Findings of the Association*  
708 *for Computational Linguistics: EMNLP 2023*, pp. 7623–7636, 2023b.
- 709 Zihao Lin, Mohammad Beigi, Hongxuan Li, Yufan Zhou, Yuxiang Zhang, Qifan Wang, Wenpeng  
710 Yin, and Lifu Huang. Navigating the dual facets: A comprehensive evaluation of sequential  
711 memory editing in large language models. *ArXiv preprint*, abs/2402.11122, 2024. URL <https://arxiv.org/abs/2402.11122>.  
712
- 713 Guofan Liu, Jinghao Zhang, Qiang Liu, Junfei Wu, Shu Wu, and Liang Wang. Uni-modal event-  
714 agnostic knowledge distillation for multimodal fake news detection. *IEEE Transactions on*  
715 *Knowledge and Data Engineering*, 2024a.  
716
- 717 Jiateng Liu, Pengfei Yu, Yuji Zhang, Sha Li, Zixuan Zhang, and Heng Ji. Evedit: Event-based  
718 knowledge editing with deductive editing boundaries. *ArXiv preprint*, abs/2402.11324, 2024b.  
719 URL <https://arxiv.org/abs/2402.11324>.
- 720 Zeyu Leo Liu, Shrey Pandit, Xi Ye, Eunsol Choi, and Greg Durrett. Codeupdatearena: Benchmarking  
721 knowledge editing on api updates. *ArXiv preprint*, abs/2407.06249, 2024c. URL <https://arxiv.org/abs/2407.06249>.  
722
- 723 Jun-Yu Ma, Jia-Chen Gu, Zhen-Hua Ling, Quan Liu, and Cong Liu. Untying the reversal curse  
724 via bidirectional language model editing. *ArXiv preprint*, abs/2310.10322, 2023. URL <https://arxiv.org/abs/2310.10322>.  
725  
726
- 727 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
728 associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.  
729
- 730 Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing  
731 memory in a transformer. In *The Eleventh International Conference on Learning Representations*,  
732 2023. URL <https://openreview.net/forum?id=MkbcAHIYgyS>.
- 733 Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. Memory-  
734 based model editing at scale. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári,  
735 Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022,*  
736 *17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning*  
737 *Research*, pp. 15817–15831. PMLR, 2022. URL <https://proceedings.mlr.press/v162/mitchell122a.html>.  
738
- 739 Qiong Nan, Qiang Sheng, Juan Cao, Yongchun Zhu, Danding Wang, Guang Yang, Jintao Li, and Kai  
740 Shu. Exploiting user comments for early detection of fake news prior to users’ commenting. *ArXiv*  
741 *preprint*, abs/2310.10429, 2023. URL <https://arxiv.org/abs/2310.10429>.  
742
- 743 Qiong Nan, Qiang Sheng, Juan Cao, Beizhe Hu, Danding Wang, and Jintao Li. Let silence speak:  
744 Enhancing fake news detection with generated comments from large language models. In *Proceed-*  
745 *ings of the 33rd ACM International Conference on Information and Knowledge Management*, pp.  
746 1732–1742, 2024.
- 747 Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. What does the knowledge neuron thesis  
748 have to do with knowledge? *ArXiv preprint*, abs/2405.02421, 2024. URL <https://arxiv.org/abs/2405.02421>.  
749
- 750 Hao Peng, Xiaozhi Wang, Chunyang Li, Kaisheng Zeng, Jiangshan Duo, Yixin Cao, Lei Hou,  
751 and Juanzi Li. Event-level knowledge editing. *ArXiv preprint*, abs/2402.13093, 2024. URL  
752 <https://arxiv.org/abs/2402.13093>.  
753
- 754 Derek Powell, Walter Gerych, and Thomas Hartvigsen. Taxi: Evaluating categorical knowledge  
755 editing for language models. *ArXiv preprint*, abs/2404.15004, 2024. URL <https://arxiv.org/abs/2404.15004>.

- 756 Siyuan Qi, Bangcheng Yang, Kailin Jiang, Xiaobo Wang, Jiaqi Li, Yifan Zhong, Yaodong Yang, and  
757 Zilong Zheng. In-context editing: Learning knowledge from self-induced distributions. *ArXiv*  
758 *preprint*, abs/2406.11194, 2024. URL <https://arxiv.org/abs/2406.11194>.  
759
- 760 Domenic Rosati, Robie Gonzales, Jinkun Chen, Xuemin Yu, Melis Erkan, Yahya Kayani,  
761 Satya Deepika Chavatapalli, Frank Rudzicz, and Hassan Sajjad. Long-form evaluation of model  
762 editing. *ArXiv preprint*, abs/2402.09394, 2024. URL <https://arxiv.org/abs/2402.09394>.
- 763 Amit Rozner, Barak Battash, Lior Wolf, and Ofir Lindenbaum. Knowledge editing in language  
764 models via adapted direct preference optimization. *arXiv preprint arXiv: 2406.09920*, 2024.  
765
- 766 Arnab Sen Sharma, David Atkinson, and David Bau. Locating and editing factual associations in  
767 mamba. *ArXiv preprint*, abs/2404.03646, 2024a. URL <https://arxiv.org/abs/2404.03646>.
- 768 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askeel, Samuel R. Bowman,  
769 Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam  
770 McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and  
771 Ethan Perez. Towards understanding sycophancy in language models. In *The Twelfth International*  
772 *Conference on Learning Representations*, 2024b. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=tvhaxkMKAn)  
773 [tvhaxkMKAn](https://openreview.net/forum?id=tvhaxkMKAn).
- 774
- 775 Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu.  
776 Retrieval-enhanced knowledge editing for multi-hop question answering in language models. *ArXiv*  
777 *preprint*, abs/2403.19631, 2024. URL <https://arxiv.org/abs/2403.19631>.
- 778 Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. Fake news detection on social media:  
779 A data mining perspective. *ACM SIGKDD explorations newsletter*, 19(1):22–36, 2017.  
780
- 781 Irene Solaiman, Zeerak Talat, William Agnew, Lama Ahmad, Dylan Baker, Su Lin Blodgett, Canyu  
782 Chen, Hal Daumé III, Jesse Dodge, Isabella Duan, et al. Evaluating the social impact of generative  
783 ai systems in systems and society. *ArXiv preprint*, abs/2306.05949, 2023. URL [https://arxiv.](https://arxiv.org/abs/2306.05949)  
784 [org/abs/2306.05949](https://arxiv.org/abs/2306.05949).
- 785 SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das.  
786 A comprehensive survey of hallucination mitigation techniques in large language models. *ArXiv*  
787 *preprint*, abs/2401.01313, 2024. URL <https://arxiv.org/abs/2401.01313>.
- 788
- 789 Rheeeya Uppaal, Apratim De, Yiting He, Yiquao Zhong, and Junjie Hu. Detox: Toxic subspace  
790 projection for model editing. *ArXiv preprint*, abs/2405.13967, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2405.13967)  
791 [abs/2405.13967](https://arxiv.org/abs/2405.13967).
- 792 Bertie Vidgen, Adarsh Agrawal, Ahmed M Ahmed, Victor Akinwande, Namir Al-Nuaimi, Najla  
793 Alfaraj, Elie Alhajar, Lora Aroyo, Trupti Bavalatti, Borhane Blili-Hamelin, et al. Introducing  
794 v0. 5 of the ai safety benchmark from mlcommons. *ArXiv preprint*, abs/2404.12241, 2024. URL  
795 <https://arxiv.org/abs/2404.12241>.
- 796
- 797 Haoran Wang, Yingdong Dou, Canyu Chen, Lichao Sun, Philip S Yu, and Kai Shu. Attacking  
798 fake news detectors via manipulating news social engagement. In *Proceedings of the ACM Web*  
799 *Conference 2023*, pp. 3978–3986, 2023a.
- 800 Haoyu Wang, Tianci Liu, Tuo Zhao, and Jing Gao. Roselora: Row and column-wise sparse low-rank  
801 adaptation of pre-trained language model for knowledge editing and fine-tuning. *ArXiv preprint*,  
802 abs/2406.10777, 2024a. URL <https://arxiv.org/abs/2406.10777>.
- 803
- 804 Jiaan Wang, Yunlong Liang, Zengkui Sun, Yuxuan Cao, and Jiarong Xu. Cross-lingual knowledge  
805 editing in large language models. *ArXiv preprint*, abs/2309.08952, 2023b. URL [https://arxiv.](https://arxiv.org/abs/2309.08952)  
806 [org/abs/2309.08952](https://arxiv.org/abs/2309.08952).
- 807 Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao, Shumin Deng, Peng Wang, Xiang Chen,  
808 Jia-Chen Gu, Yong Jiang, Pengjun Xie, et al. Knowledge mechanisms in large language models: A  
809 survey and perspective. *ArXiv preprint*, abs/2407.15017, 2024b. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2407.15017)  
[2407.15017](https://arxiv.org/abs/2407.15017).

- 810 Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang,  
811 Linyi Yang, Jindong Wang, and Huajun Chen. Detoxifying large language models via knowledge  
812 editing. *ArXiv preprint*, abs/2403.14472, 2024c. URL <https://arxiv.org/abs/2403.14472>.  
813
- 814 Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang,  
815 and Huajun Chen. Wise: Rethinking the knowledge memory for lifelong model editing of large  
816 language models. *ArXiv preprint*, abs/2405.14768, 2024d. URL <https://arxiv.org/abs/2405.14768>.  
817
- 818 Renzhi Wang and Piji Li. Lemoe: Advanced mixture of experts adaptor for lifelong model editing of  
819 large language models. *ArXiv preprint*, abs/2406.20030, 2024a. URL <https://arxiv.org/abs/2406.20030>.  
820
- 821 Renzhi Wang and Piji Li. Semantic are beacons: A semantic perspective for unveiling parameter-  
822 efficient fine-tuning in knowledge learning. *ArXiv preprint*, abs/2405.18292, 2024b. URL <https://arxiv.org/abs/2405.18292>.  
823
- 824 Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, et al. Knowledge editing  
825 for large language models: A survey. *ArXiv preprint*, abs/2310.16218, 2023c. URL <https://arxiv.org/abs/2310.16218>.  
826
- 827 Wenxuan Wang, Juluan Shi, Zhaopeng Tu, Youliang Yuan, Jen-tse Huang, Wenxiang Jiao, and  
828 Michael R Lyu. The earth is flat? unveiling factual errors in large language models. *ArXiv preprint*,  
829 abs/2401.00761, 2024e. URL <https://arxiv.org/abs/2401.00761>.  
830
- 831 Xiaohan Wang, Shengyu Mao, Ningyu Zhang, Shumin Deng, Yunzhi Yao, Yue Shen, Lei Liang,  
832 Jinjie Gu, and Huajun Chen. Editing conceptual knowledge for large language models. *ArXiv*  
833 *preprint*, abs/2403.06259, 2024f. URL <https://arxiv.org/abs/2403.06259>.  
834
- 835 Yiwei Wang, Muhao Chen, Nanyun Peng, and Kai-Wei Chang. Deepedit: Knowledge editing as  
836 decoding with constraints. *ArXiv preprint*, abs/2401.10471, 2024g. URL <https://arxiv.org/abs/2401.10471>.  
837
- 838 Yifan Wei, Xiaoyan Yu, Huanhuan Ma, Fangyu Lei, Yixuan Weng, Ran Song, and Kang Liu. Assess-  
839 ing knowledge editing in language models via relation perspective. *ArXiv preprint*, abs/2311.09053,  
840 2023. URL <https://arxiv.org/abs/2311.09053>.  
841
- 842 Zihao Wei, Jingcheng Deng, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. Mlake: Mul-  
843 tilingual knowledge editing benchmark for large language models. *ArXiv preprint*, abs/2404.04990,  
844 2024a. URL <https://arxiv.org/abs/2404.04990>.  
845
- 846 Zihao Wei, Liang Pang, Hanxing Ding, Jingcheng Deng, Huawei Shen, and Xueqi Cheng. Stable  
847 knowledge editing in large language models. *ArXiv preprint*, abs/2402.13048, 2024b. URL  
848 <https://arxiv.org/abs/2402.13048>.  
849
- 850 Suhang Wu, Minlong Peng, Yue Chen, Jinsong Su, and Mingming Sun. Eva-kellm: A new benchmark  
851 for evaluating knowledge editing of llms. *ArXiv preprint*, abs/2308.09954, 2023. URL <https://arxiv.org/abs/2308.09954>.  
852
- 853 Xiaobao Wu, Liangming Pan, William Yang Wang, and Anh Tuan Luu. Updating language models  
854 with unstructured facts: Towards practical knowledge editing. *ArXiv preprint*, abs/2402.18909,  
855 2024. URL <https://arxiv.org/abs/2402.18909>.  
856
- 857 Jiakuan Xie, Pengfei Cao, Yuheng Chen, Yubo Chen, Kang Liu, and Jun Zhao. Memla: Enhancing  
858 multilingual knowledge editing with neuron-masked low-rank adaptation. *arXiv preprint arXiv:*  
859 *2406.11566*, 2024.  
860
- 861 Derong Xu, Ziheng Zhang, Zhihong Zhu, Zhenxi Lin, Qidong Liu, Xian Wu, Tong Xu, Xiangyu  
862 Zhao, Yefeng Zheng, and Enhong Chen. Editing factual knowledge and explanatory ability of  
863 medical large language models. *ArXiv preprint*, abs/2402.18099, 2024. URL <https://arxiv.org/abs/2402.18099>.



- 864 Jianhao Yan, Futing Wang, Yafu Li, and Yue Zhang. Potential and challenges of model editing for  
865 social debiasing. *ArXiv preprint*, abs/2402.13462, 2024. URL [https://arxiv.org/abs/2402.](https://arxiv.org/abs/2402.13462)  
866 [13462](https://arxiv.org/abs/2402.13462).
- 867 Wanli Yang, Fei Sun, Xinyu Ma, Xun Liu, Dawei Yin, and Xueqi Cheng. The butterfly effect of model  
868 editing: Few edits can trigger large language models collapse. *ArXiv preprint*, abs/2402.09656,  
869 2024. URL <https://arxiv.org/abs/2402.09656>.
- 870 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen,  
871 and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. *ArXiv*  
872 *preprint*, abs/2305.13172, 2023. URL <https://arxiv.org/abs/2305.13172>.
- 873 Yunzhi Yao, Ningyu Zhang, Zekun Xi, Mengru Wang, Ziwen Xu, Shumin Deng, and Huajun Chen.  
874 Knowledge circuits in pretrained transformers. *ArXiv preprint*, abs/2405.17969, 2024. URL  
875 <https://arxiv.org/abs/2405.17969>.
- 876 Xunjian Yin, Jin Jiang, Liming Yang, and Xiaojun Wan. History matters: Temporal knowledge  
877 editing in large language model. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
878 volume 38, pp. 19413–19421, 2024.
- 879 Lang Yu, Qin Chen, Jie Zhou, and Liang He. Melo: Enhancing model editing with neuron-indexed  
880 dynamic lora. *ArXiv preprint*, abs/2312.11795, 2023. URL [https://arxiv.org/abs/2312.](https://arxiv.org/abs/2312.11795)  
881 [11795](https://arxiv.org/abs/2312.11795).
- 882 Zhenrui Yue, Huimin Zeng, Yimeng Lu, Lanyu Shang, Yang Zhang, and Dong Wang. Evidence-  
883 driven retrieval augmented response generation for online misinformation. In *Proceedings of the*  
884 *2024 Conference of the North American Chapter of the Association for Computational Linguistics:*  
885 *Human Language Technologies (Volume 1: Long Papers)*, pp. 5628–5643, 2024.
- 886 Mengqi Zhang, Bowen Fang, Qiang Liu, Pengjie Ren, Shu Wu, Zhumin Chen, and Liang Wang.  
887 Enhancing multi-hop reasoning through knowledge erasure in large language model editing. *ArXiv*  
888 *preprint*, abs/2408.12456, 2024a. URL <https://arxiv.org/abs/2408.12456>.
- 889 Mengqi Zhang, Xiaotian Ye, Qiang Liu, Pengjie Ren, Shu Wu, and Zhumin Chen. Knowledge  
890 graph enhanced large language model editing. *ArXiv preprint*, abs/2402.13593, 2024b. URL  
891 <https://arxiv.org/abs/2402.13593>.
- 892 Mengqi Zhang, Xiaotian Ye, Qiang Liu, Pengjie Ren, Shu Wu, and Zhumin Chen. Uncovering  
893 overfitting in large language model editing. *ArXiv preprint*, abs/2410.07819, 2024c. URL <https://arxiv.org/abs/2410.07819>.
- 894 Ningyu Zhang, Zekun Xi, Yujie Luo, Peng Wang, Bozhong Tian, Yunzhi Yao, Jintian Zhang, Shumin  
895 Deng, Mengshu Sun, Lei Liang, et al. Oneedit: A neural-symbolic collaboratively knowledge  
896 editing system. *ArXiv preprint*, abs/2409.07497, 2024d. URL [https://arxiv.org/abs/2409.](https://arxiv.org/abs/2409.07497)  
897 [07497](https://arxiv.org/abs/2409.07497).
- 898 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi,  
899 Shengyu Mao, Jintian Zhang, Yuansheng Ni, et al. A comprehensive study of knowledge editing  
900 for large language models. *ArXiv preprint*, abs/2401.01286, 2024e. URL [https://arxiv.org/](https://arxiv.org/abs/2401.01286)  
901 [abs/2401.01286](https://arxiv.org/abs/2401.01286).
- 902 Shaolei Zhang, Tian Yu, and Yang Feng. Truthx: Alleviating hallucinations by editing large language  
903 models in truthful space. *ArXiv preprint*, abs/2402.17811, 2024f. URL [https://arxiv.org/](https://arxiv.org/abs/2402.17811)  
904 [abs/2402.17811](https://arxiv.org/abs/2402.17811).
- 905 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,  
906 Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi.  
907 Siren’s song in the ai ocean: A survey on hallucination in large language models. *ArXiv preprint*,  
908 abs/2309.01219, 2023. URL <https://arxiv.org/abs/2309.01219>.
- 909 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,  
910 Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen,  
911 Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and  
912 Ji-Rong Wen. A survey of large language models. *ArXiv preprint*, abs/2303.18223, 2023. URL  
913 <https://arxiv.org/abs/2303.18223>.

918 Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can  
919 we edit factual knowledge by in-context learning? In Houda Bouamor, Juan Pino, and Kalika Bali  
920 (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*,  
921 pp. 4862–4876, Singapore, 2023. Association for Computational Linguistics. doi: 10.18653/v1/  
922 2023.emnlp-main.296. URL <https://aclanthology.org/2023.emnlp-main.296>.

923 Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen.  
924 Mquake: Assessing knowledge editing in language models via multi-hop questions. *ArXiv*  
925 *preprint*, abs/2305.14795, 2023. URL <https://arxiv.org/abs/2305.14795>.

926  
927 Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv  
928 Kumar. Modifying memories in transformer models. *ArXiv preprint*, abs/2012.00363, 2020. URL  
929 <https://arxiv.org/abs/2012.00363>.

930 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,  
931 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering:  
932 A top-down approach to ai transparency. *ArXiv preprint*, abs/2310.01405, 2023. URL <https://arxiv.org/abs/2310.01405>.  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

# Content of Appendix

- A Reproducibility Statement** **20**
  
- B Details of the Benchmarked Knowledge Editing Techniques** **21**
  
- C A More Detailed Related Work** **22**
  
- D Impact Statement** **22**
  
- E More Experiment Results** **23**
  - E.1 Generalization Scores of Knowledge Editing Methods on Each Domain . . . . . **23**
  - E.2 Portability Scores of Knowledge Editing Methods on More Domains . . . . . **28**
  - E.3 Robustness Scores of Knowledge Editing Methods on More Domains . . . . . **31**
  
- F Examples of HalluEditBench** **34**

## A REPRODUCIBILITY STATEMENT

We conduct the experiments on NVIDIA RTX A6000 GPUs. The decoding temperatures are 0 to ensure reproducibility. The model checkpoints are downloaded from <https://huggingface.co/>. The specific download links are as follows:

- Llama2-7B: <https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>
- Llama3-8B: <https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>
- Mistral-v0.3-7B: <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

We adopt GPT-4o with the prompt below to generate *Generalization* and *Locality* evaluation questions:

---

Given a fact triplet (subject, relation, object), a question asking for the object, and a wrong answer, the correct answer to the question should be the object in the triplet.

Generate the following types of questions:

1. Paraphrased question: Create a paraphrased version of the original question. The correct answer should still be the object from the triplet.
2. Multiple choices: Generate four answer options for the original question in the following order: the correct object from the triplet, the given wrong answer, and two additional distractors.
3. Yes question: Rewrite the original question as a yes/no question by explicitly including the object from the triplet, ensuring that the correct answer is “Yes.”
4. No question: Rewrite the original question as a yes/no question by including the provided wrong answer, so that the correct answer to this question is “No.”
5. Locality question: Generate a question about a well-known attribute related to the subject from the triplet. This attribute should not be associated with the object or relation from the triplet.
6. Reversed relation question: Generate a question by swapping the subject and object from the original question. The answer should now be the subject from the triplet.

Output the result in JSON format with the following keys: “paraphrased\_question”, “multiple\_choices”, “yes\_question”, “no\_question”, “locality\_question”, and “reversed\_relation\_question.”

---

We adopt GPT-4o with the following prompt to generate evaluation questions in *Portability* aspect.

---

Given a subject, a relation, a 1-hop question, and its answer, create 2-hop, 3-hop, 4-hop, 5-hop, and 6-hop questions, along with their correct answers.

Always use the provided subject and relation to create multi-hop questions and include the preceding question in the subsequent question (for example, include the 2-hop question in 3-hop question, include the 3-hop question in 4-hop question).

DO NOT include the correct answer to any previous multi-hop question in subsequent ones (for example, do not include the correct answer to the 2-hop question in the 3-hop or 4-hop questions).

Ensure that the answers for all multi-hop questions are accurate, and do not use ‘N/A’ as an answer.

You must include the given subject and relation in all of the 2-hop, 3-hop, 4-hop, 5-hop, and 6-hop questions. Output in JSON format. An example is provided below:

Example input:

subject: Amazon, relation: founder

1hop\_question: Who is the Amazon founder? 1hop\_answer: Jeff Bezos

Example output:

```
{
  "2hop_question": "Who is the spouse of the Amazon founder?",
  "2hop_answer": "MacKenzie Scott",
  "3hop_question": "Which university did the spouse of the Amazon founder attend for their undergraduate studies?",
  "3hop_answer": "Princeton University",
  "4hop_question": "In which city is the university that the spouse of the Amazon founder attended located?",
  "4hop_answer": "Princeton",
  "5hop_question": "In which state is the city located where the university that the spouse of the Amazon founder attended is situated?",
  "5hop_answer": "New Jersey",
  "6hop_question": "In which country is the state located where the city is situated that contains the university the spouse of the Amazon founder attended?",
  "6hop_answer": "United States",
}
```

---

## B DETAILS OF THE BENCHMARKED KNOWLEDGE EDITING TECHNIQUES

**FT-L** (Zhu et al., 2020; Meng et al., 2022) Constrained Fine-Tuning (FT-L) is a targeted approach to fine-tuning that focuses on adjusting a specific layer within a model’s feed-forward network (FFN). Guided by causal tracing results from ROME, FT-L modifies the layer most associated with the desired changes. The goal of FT-L is to fine-tune the model by maximizing the likelihood of the target sequence, particularly focusing on the prediction of the last token, ensuring that the model adapts to modified facts without affecting its broader performance. To achieve this, explicit parameter-space norm constraints are applied to the weights, ensuring minimal interference with unmodified facts and preserving the integrity of the model’s original knowledge.

**FT-M** (Zhang et al., 2024e) In contrast to FT-L, which fine-tunes by maximizing the probability of all tokens in the target sequence based on the last token’s prediction, Fine-Tuning with Masking (FT-M) refines this approach to align more closely with the traditional fine-tuning objective. FT-M also targets the same FFN layer identified by causal tracing but employs a masked training strategy. Specifically, it uses cross-entropy loss on the target answer while masking out the original text, ensuring that the model is trained directly on the relevant target content. This approach mitigates potential deviations from the original fine-tuning objective and provides a more precise adjustment of the model’s weights with minimal disruption to unrelated model behavior.

**LoRA** (Hu et al., 2022) Low-Rank Adaptation (LoRA) is a parameter-efficient fine-tuning method that enhances training efficiency by introducing trainable rank decomposition matrices into Transformer layers. Rather than updating the original model parameters directly, LoRA focuses on training expansion and reduction matrices with low intrinsic rank, which allows for significant dimensionality reduction and thus faster training. Specifically, LoRA freezes the pretrained model weights and optimizes rank decomposition matrices to indirectly adapt dense layers without altering the original parameters. This approach greatly reduces the number of trainable parameters needed for downstream tasks, enabling more efficient training and lowering hardware requirements.

**ROME** (Meng et al., 2022) Rank-One Model Editing (ROME) is a “Locate-then-Edit” technique designed to modify factual associations within transformer models. ROME localizes these associations along three key dimensions: (1) the MLP module parameters, (2) within a range of middle layers, and (3) specifically during the processing of the last token of the subject. It employs causal intervention to trace the causal effects of hidden state activations, identifying the specific modules that mediate the recall of factual information. Once these decisive MLP modules are localized, ROME makes small, targeted rank-one changes to the parameters of a single MLP module, effectively altering individual factual associations while minimizing disruption to the overall model behavior. This precise parameter adjustment enables direct updates to the model’s factual knowledge.

**MEMIT** (Meng et al., 2023) Mass Editing Memory in a Transformer (MEMIT) builds upon ROME to generalize the editing of feedforward networks (FFNs) in pre-trained transformer models for mass knowledge updates. While ROME focuses on localizing and modifying factual associations within single layers, MEMIT extends this strategy to perform mass edits across a range of critical layers. MEMIT uses causal tracing to identify MLP layers that act as mediators of factual recall, similarly to ROME, but scales the process to enable the simultaneous insertion of thousands of new memories. By explicitly calculating parameter updates, MEMIT targets these critical layers and updates them efficiently, offering a scalable multi-layer update algorithm that enhances and expands upon ROME’s capability to modify knowledge across many memories concurrently, achieving orders of magnitude greater scalability.

**ICE** (Zheng et al., 2023) In-Context Knowledge Editing (IKE) leverages in-context learning (ICL) to modify model outputs without altering the model’s parameters. This approach reduces computational overhead and avoids potential side effects from parameter updates, offering a more efficient and safer way to modify knowledge in large language models. IKE enhances interpretability, providing a human-understandable method for calibrating model behaviors. It achieves this by constructing three types of demonstrations-copy, update, and retain-that guide the model in producing reliable fact editing through the use of a demonstration store. This store, built from training examples, allows the model to retrieve the most relevant demonstrations to inform its responses, improving accuracy in modifying specific factual outputs. In-Context Editing (ICE) is a simple baseline variant of IKE, which directly uses the new fact as context without additional demonstrations.

1134 **GRACE** (Hartvigsen et al., 2024) GRACE is a knowledge editing method designed to enable  
1135 thousands of sequential edits without the pitfalls of overfitting or loss of previously learned knowledge,  
1136 which are common in conventional knowledge editing approaches. GRACE introduces an adaptor to  
1137 a chosen layer of a model, allowing for layer-to-layer transformation adjustments without altering the  
1138 model’s original weights. This adaptor caches embeddings corresponding to input errors and learns  
1139 values that map to the desired model outputs, effectively functioning as a codebook where edits are  
1140 stored. The codebook of edits maintains model stability and allows for more extended sequences of  
1141 edits. GRACE includes a deferral mechanism that decides whether to use the codebook for a given  
1142 input, enabling the model to dynamically search and replace hidden states based on stored knowledge.  
1143 This approach allows for flexible and efficient updates to the models predictions while preserving its  
1144 pre-trained capabilities.

## 1145 C A MORE DETAILED RELATED WORK

1148 Knowledge Editing has been adopted as one of the mainstream paradigms to address the hallucinations  
1149 in LLMs efficiently (Chen & Shu, 2024a; Tonmoy et al., 2024). Besides benchmarks, recent works  
1150 have studied knowledge editing from different perspectives. The first line of works aims to probe into  
1151 the relationship between localization and editing and gain a deeper understanding of the working  
1152 mechanisms of different techniques (Wang et al., 2024b; Niu et al., 2024; Hase et al., 2024a;b;  
1153 Ferrando et al., 2024; Gupta et al., 2024; Chen et al., 2024d;c; Zou et al., 2023; Yao et al., 2024).  
1154 For example, Hase et al. (2024a) found that *Causal Tracing* actually does not provide any insight  
1155 into which MLP layer is the best option to edit. The second line of works intends to enhance the  
1156 performance and applicability of knowledge editing in specific scenarios (Rozner et al., 2024; Jiang  
1157 et al., 2024a;b; Zhang et al., 2024c;b;d;a;f; Wu et al., 2024; Qi et al., 2024; Sharma et al., 2024a;  
1158 Li et al., 2024c;b; Fang et al., 2024; Wang & Li, 2024a;b; Wang et al., 2024g;f;d; 2023b; Cheng  
1159 et al., 2024b;a; Xie et al., 2024; Bi et al., 2024c;b;a; Chen et al., 2024b; Wei et al., 2024b; Fei et al.,  
1160 2024; Xu et al., 2024; Gu et al., 2023; Yin et al., 2024; Cai et al., 2024a; Liu et al., 2024b; Ge  
1161 et al., 2024b; Deng et al., 2024; Peng et al., 2024). For example, Ma et al. (2023) proposed a new  
1162 method named Bidirectionally Inversible Relationship Modeling (BIRD) to mitigate the *reversal*  
1163 *curse* issue in bidirectional language model editing and improve the performance. The third line  
1164 of works investigates the side effect of knowledge editing techniques (Hsueh et al., 2024; Gu et al.,  
1165 2024; Hoelscher-Obermaier et al., 2023; Hua et al., 2024; Yang et al., 2024; Li et al., 2023a; Cohen  
1166 et al., 2024). For example, Yang et al. (2024) discovered that even one single edit could cause a  
1167 significant performance degradation in mainstream benchmarks. The fourth line of works explores the  
1168 potential misuse risks of knowledge editing or its applications beyond correcting hallucinations (Chen  
1169 et al., 2024a; Uppaal et al., 2024; Wang et al., 2024c; Cai et al., 2024b; Yan et al., 2024). For  
1170 example, Chen et al. (2024a) proposed to reformulate knowledge editing as a new type of safety  
1171 threat, namely *Editing Attack*, and validated its risk of injecting misinformation or bias into LLMs  
1172 stealthily, suggesting the feasibility of disseminating misinformation or bias with LLMs as new  
1173 channels. The social impact of knowledge editing techniques, especially on safety aspect, is worth  
1174 more attention (Solaiman et al., 2023; Vidgen et al., 2024).

## 1175 D IMPACT STATEMENT

1176 Misinformation is a longstanding threat for online safety and public trust (Chen et al., 2022; Wang  
1177 et al., 2023a). The conventional countermeasures include *detection* (Shu et al., 2017; Nan et al.,  
1178 2024; 2023; Liu et al., 2024a), *intervention* (Bak-Coleman et al., 2022; Aghajari et al., 2023;  
1179 Hartwig et al., 2024; Yue et al., 2024; He et al., 2023) and *attribution* (Huang et al., 2024a;b; Beigi  
1180 et al., 2024). Hallucinations, which could be defined as the non-factual information unintentionally  
1181 generated by LLMs when used by normal users (Chen & Shu, 2024a;b), have become a new type of  
1182 misinformation and may cause severe information pollution to the online space. Knowledge editing is  
1183 a promising paradigm to correct hallucinations and contribute to the fight against the misinformation  
1184 crisis in the era of LLMs, due to its advantage of avoiding retraining from scratch. However, our  
1185 work sheds light on the potential limitations of current knowledge editing techniques and calls for  
1186 more effort to address these challenges collectively in the future.

## E MORE EXPERIMENT RESULTS

### E.1 GENERALIZATION SCORES OF KNOWLEDGE EDITING METHODS ON EACH DOMAIN

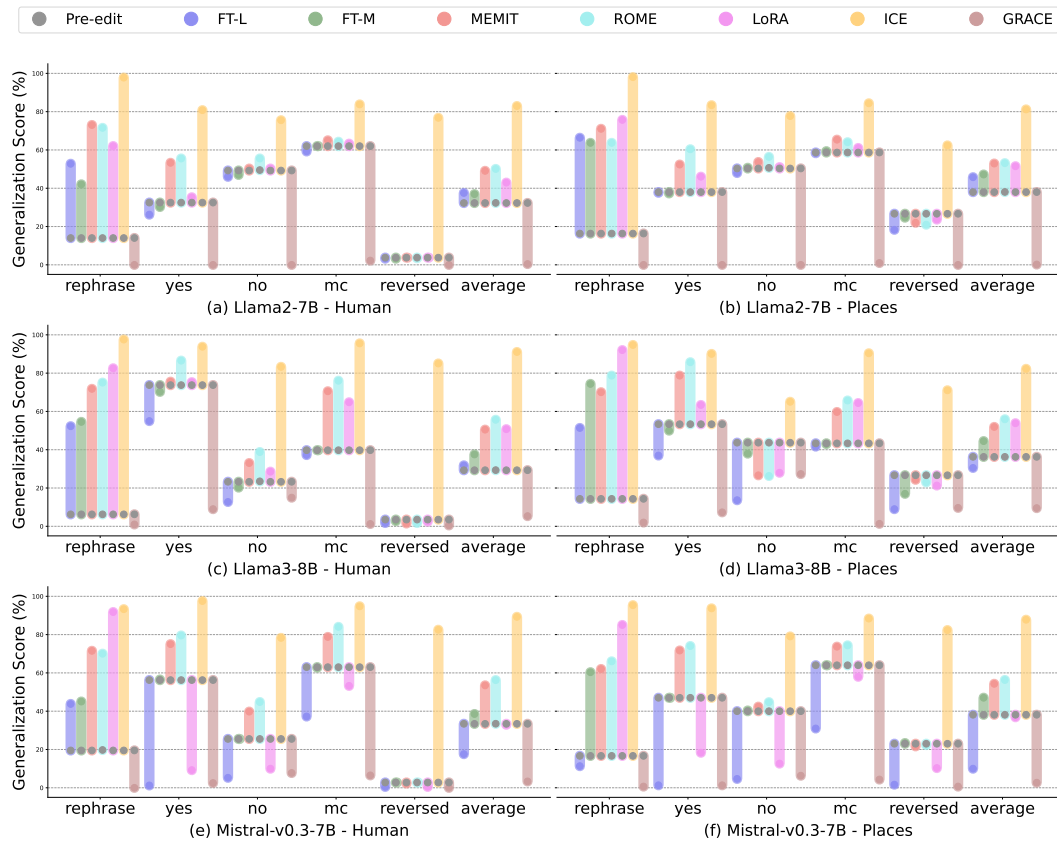


Figure 8: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), Reversed Questions (“reversed”). The “average” refers to the averaged scores over five types of questions. The domains include “human” and “places”.

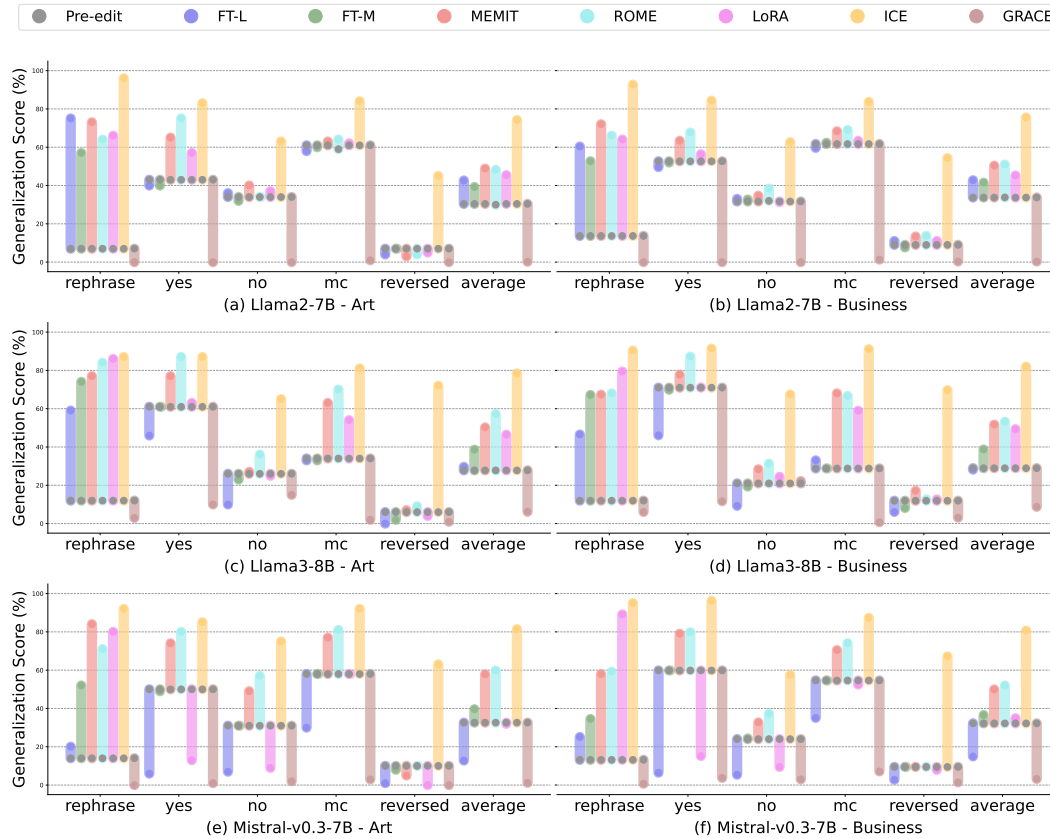


Figure 9: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), Reversed Questions (“reversed”). The “average” refers to the averaged scores over five types of questions. The domains include “art” and “business”.



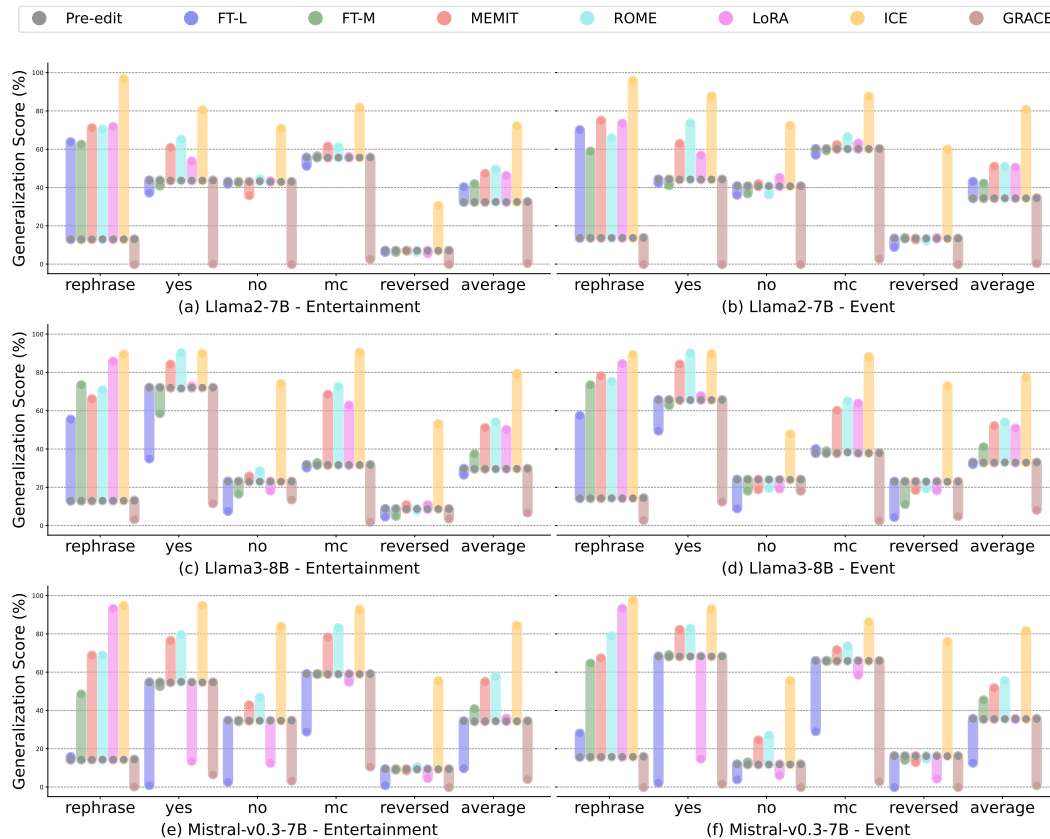


Figure 10: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), Reversed Questions (“reversed”). The “average” refers to the averaged scores over five types of questions. The domains include “entertainment” and “event”.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

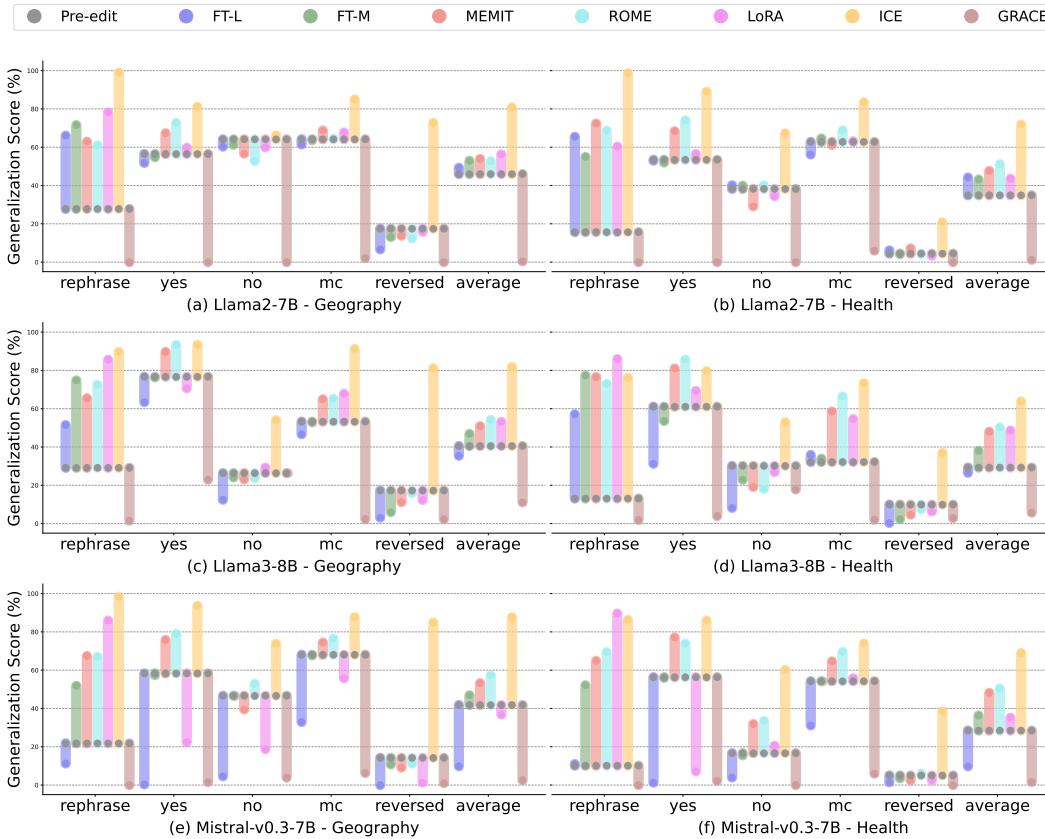


Figure 11: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), Reversed Questions (“reversed”). The “average” refers to the averaged scores over five types of questions. The domains include “geography” and “health”.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

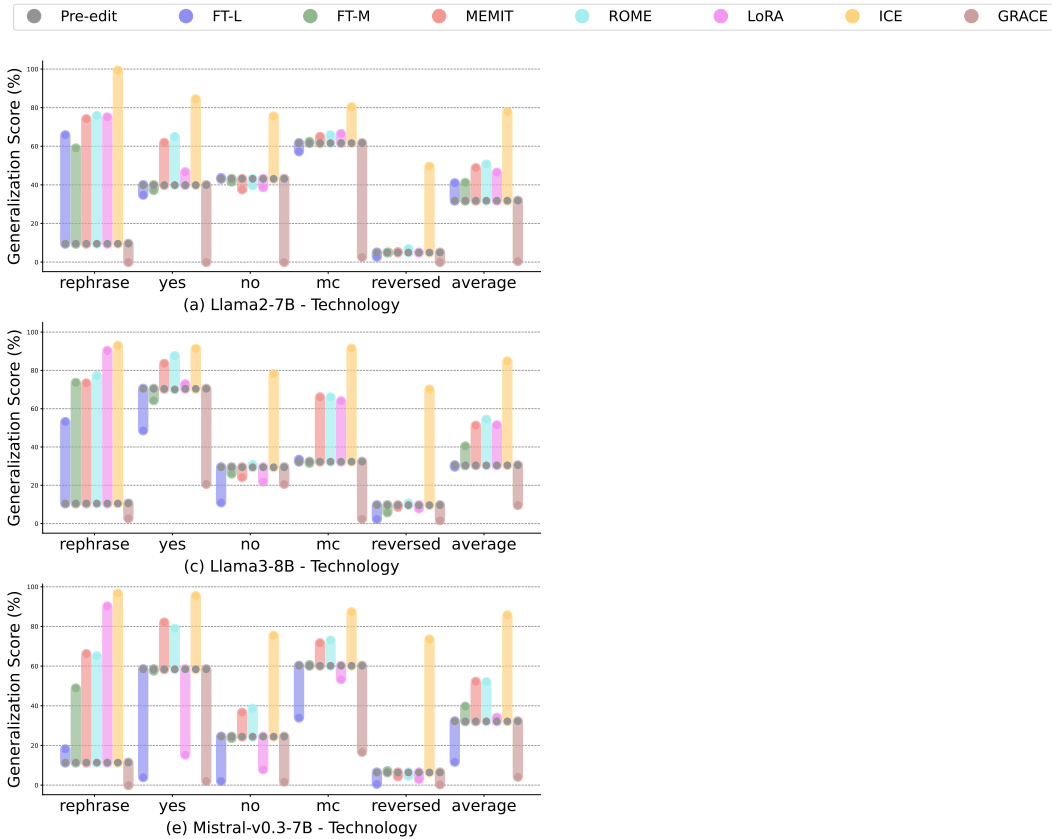


Figure 12: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), Reversed Questions (“reversed”). The “average” refers to the averaged scores over five types of questions. The domain is “technology”.

## E.2 PORTABILITY SCORES OF KNOWLEDGE EDITING METHODS ON MORE DOMAINS

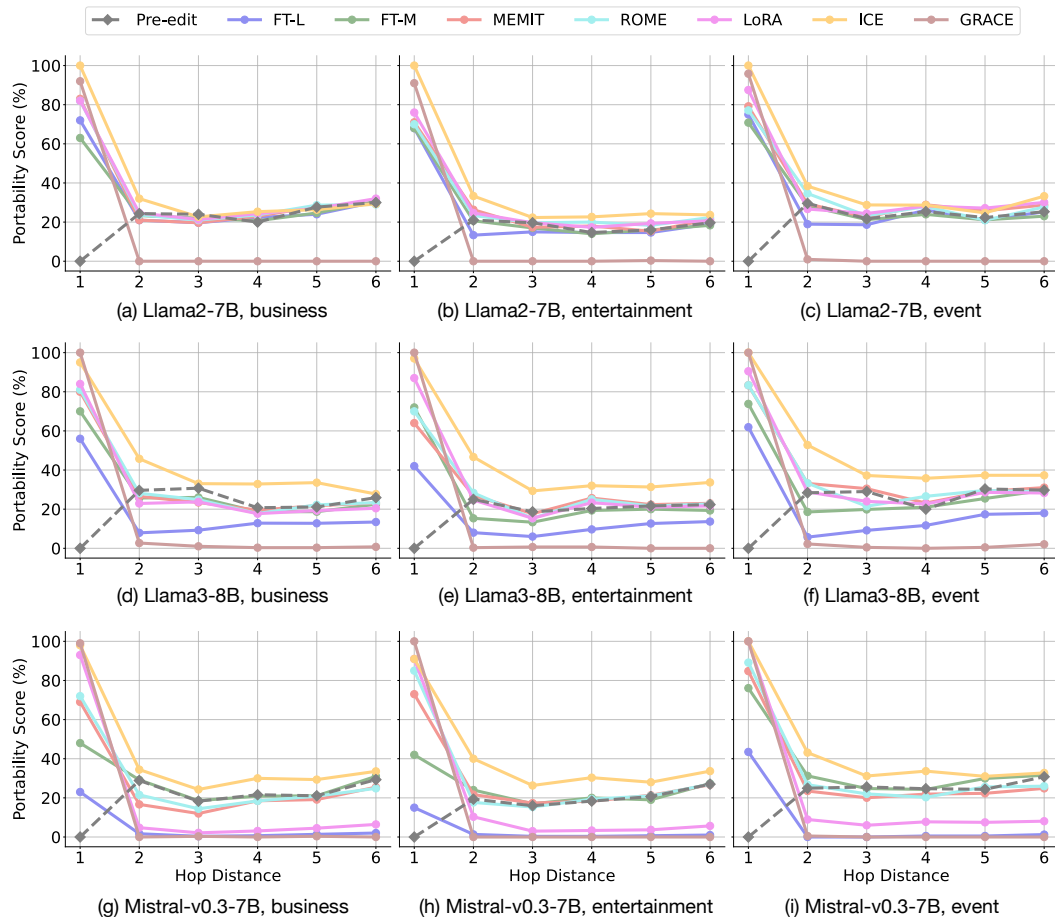


Figure 13: **Portability Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions when with  $N$  hops. The Portability Evaluation Questions are the same as Efficacy Evaluation Questions when  $N$  is 1. The domains include “business”, “entertainment”, and “event”.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

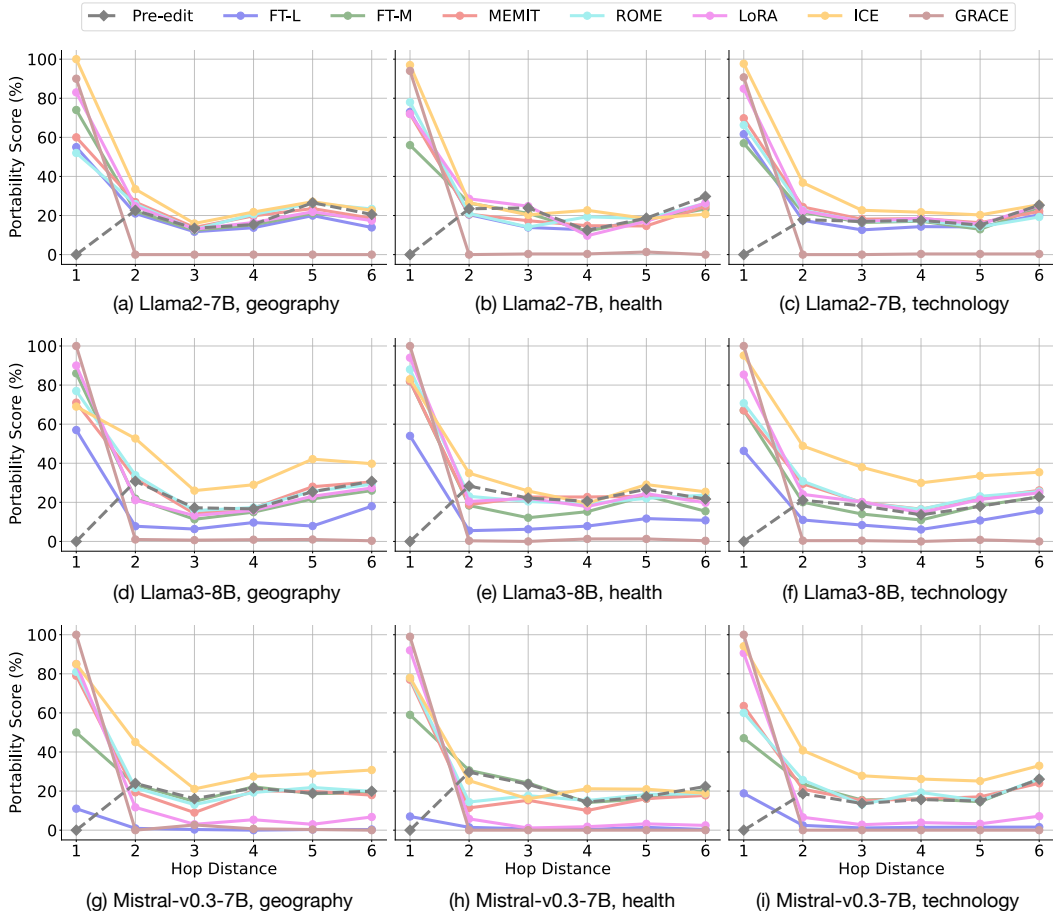


Figure 14: **Portability Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions when with  $N$  hops. The Portability Evaluation Questions are the same as Efficacy Evaluation Questions when  $N$  is 1. The domains include “geography”, “health”, and “technology”.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

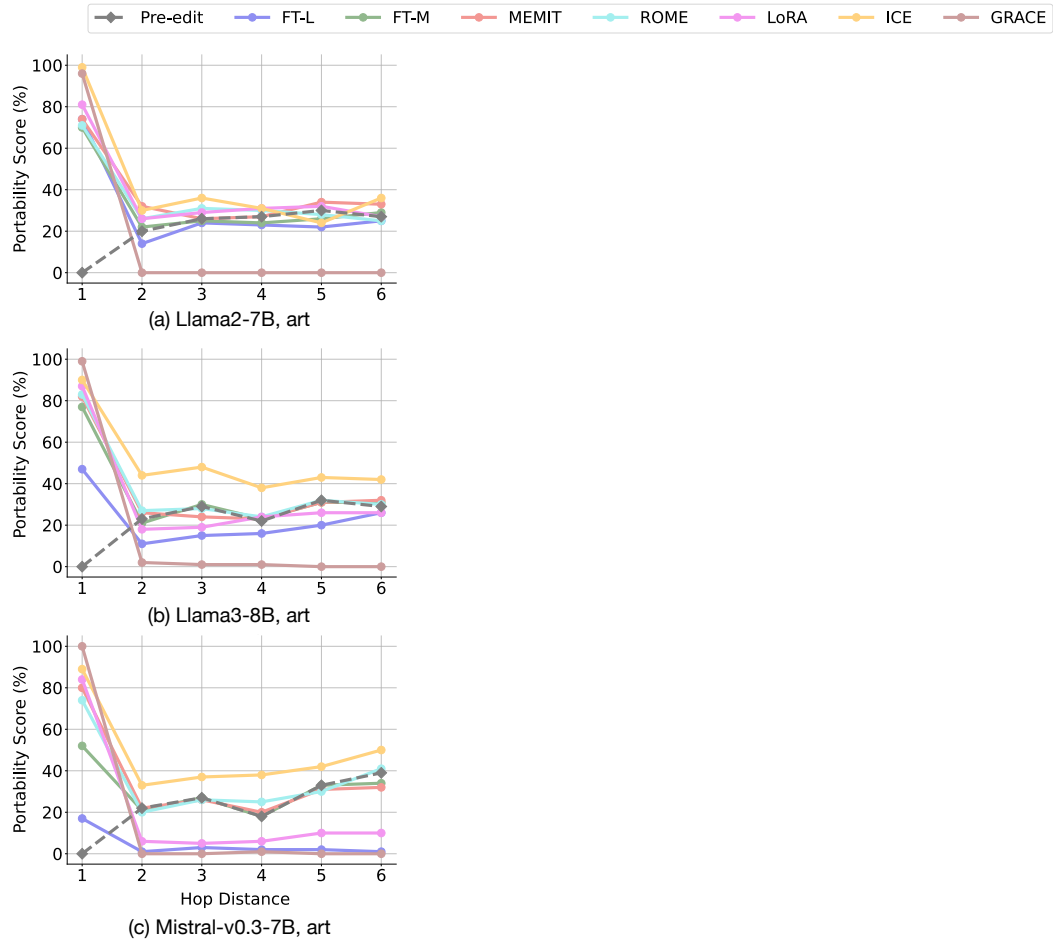


Figure 15: **Portability Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions when with  $N$  hops. The Portability Evaluation Questions are the same as Efficacy Evaluation Questions when  $N$  is 1. The domain is “art”.

## E.3 ROBUSTNESS SCORES OF KNOWLEDGE EDITING METHODS ON MORE DOMAINS

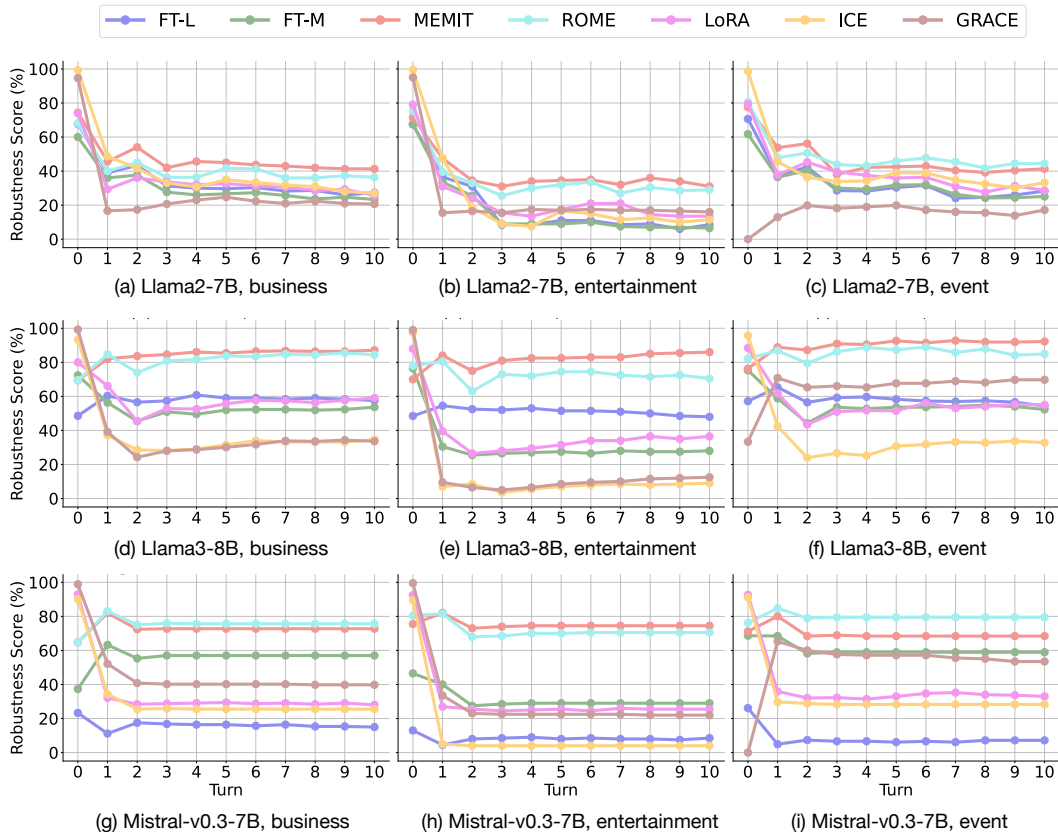


Figure 16: **Robustness Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with  $M$  turns ( $M = 1 \sim 10$ ). We regard Efficacy Scores as the Robustness Scores when  $M$  is 0. The domains include “business”, “entertainment”, and “event”.

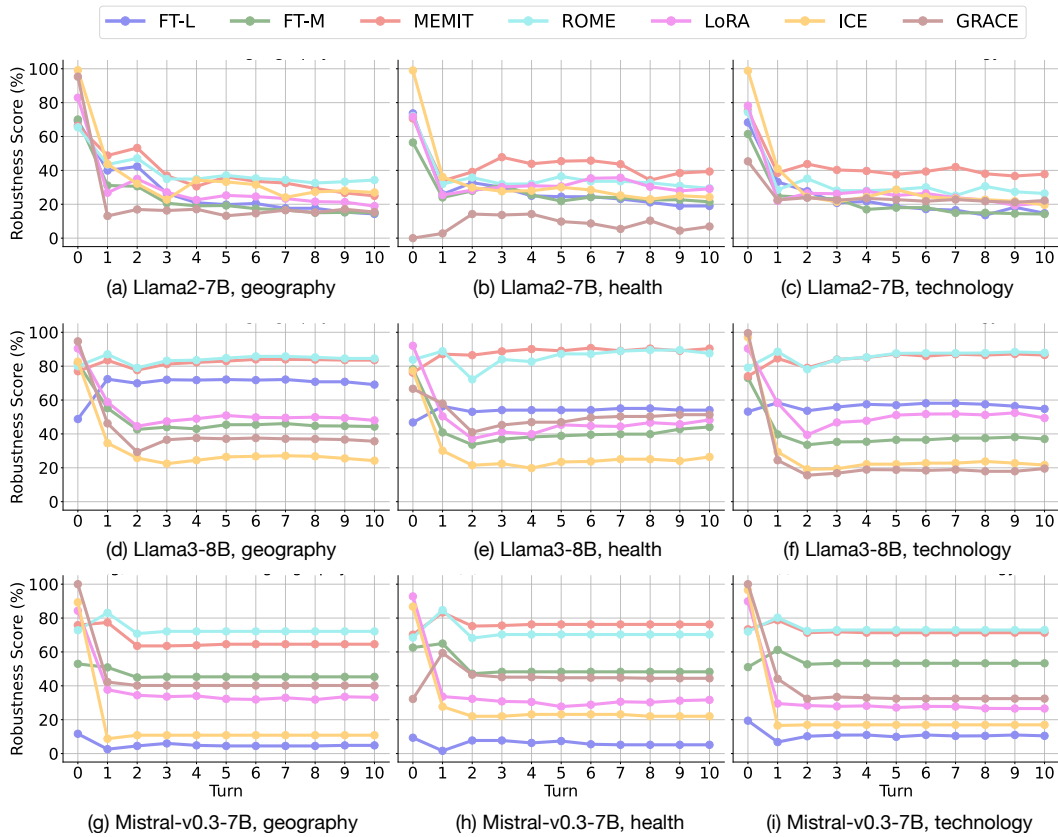


Figure 17: **Robustness Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with  $M$  turns ( $M = 1 \sim 10$ ). We regard Efficacy Scores as the Robustness Scores when  $M$  is 0. The domains include “geography”, “health”, and “technology”.



1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

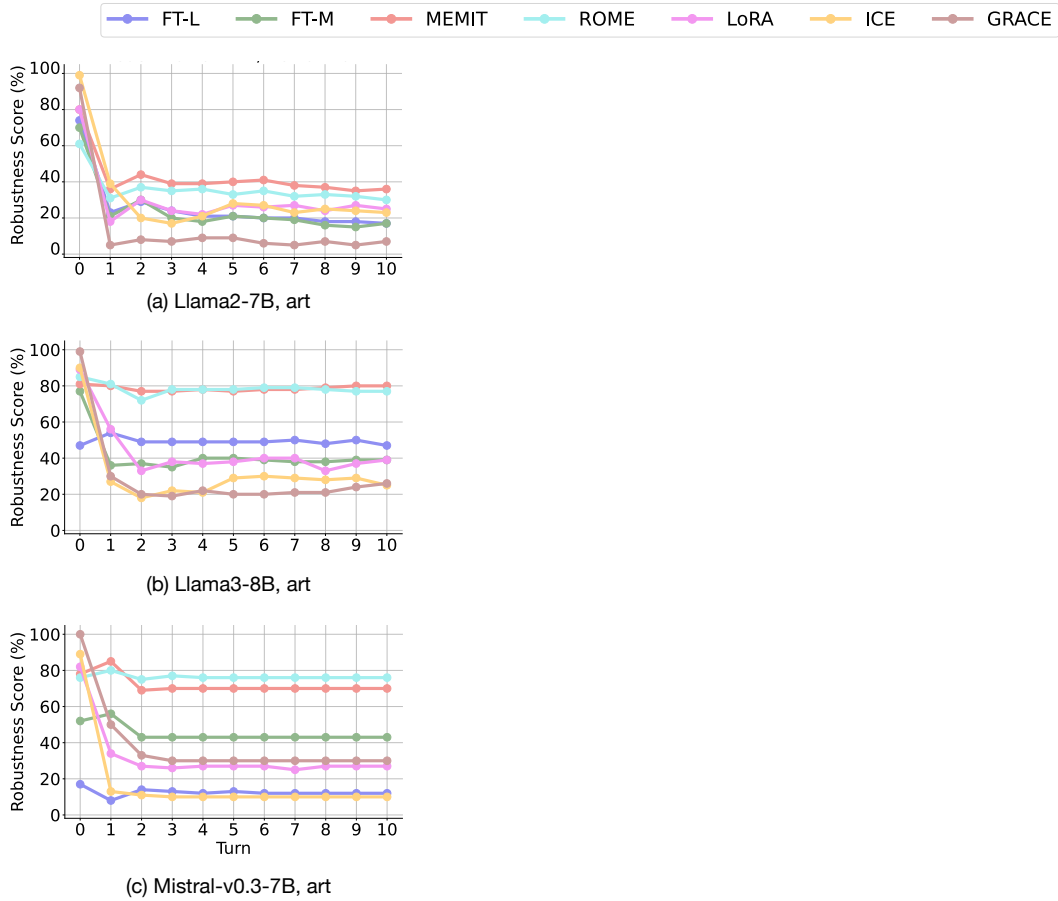


Figure 18: **Robustness Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with  $M$  turns ( $M = 1 \sim 10$ ). We regard Efficacy Scores as the Robustness Scores when  $M$  is 0. The domain is “art”.

1782 F EXAMPLES OF HalluEditBench  
1783

1784 The examples of evaluation questions based on Llama3-8B’s hallucinations are shown as follows:  
1785

---

1786 Domain: Places, Topic: Country, Knowledge Triplet: (Trinidad and Tobago, highest point, El Cerro  
1787 del Aripo)  
1788

1789 **Efficacy Evaluation Questions:** What is the highest point of Trinidad and Tobago?  
1790 Ground-truth Answer: El Cerro del Aripo  
1791 Hallucinated Answer of Llama3-8B before editing: Pierrepont Hill  
1792

1793 **Generalization Evaluation Questions:**  
1794 Rephrased Question: What is the name of the highest peak in Trinidad and Tobago?  
1795 Ground-truth Answer: El Cerro del Aripo  
1796

1797 Yes-or-No Questions with Yes as answer: Is El Cerro del Aripo the highest point of Trinidad and  
1798 Tobago?  
1799 Ground-truth Answer: Yes  
1800

1801 Yes-or-No Questions with No as answer: Is Pierrepont Hill the highest point of Trinidad and Tobago?  
1802 Ground-truth Answer: No  
1803

1804 Multi-Choice Questions: What is the highest point of Trinidad and Tobago? A. Mount Kilimanjaro B.  
1805 Mount Everest C. El Cerro del Aripo D. Pierrepont Hill  
1806 Ground-truth Answer: C  
1807

1808 Reversed Questions: Which country has El Cerro del Aripo as its highest point?  
1809 Ground-truth Answer: Trinidad and Tobago  
1810

1811 **Portability Evaluation Questions:**  
1812 One-hop question: What is the highest point of Trinidad and Tobago?  
1813 Ground-truth Answer: El Cerro del Aripo  
1814

1815 Two-hop question: In which mountain range is the highest point of Trinidad and Tobago located?  
1816 Ground-truth Answer: Northern Range  
1817

1818 Three-hop question: Which island is the mountain range located in that contains the highest point of  
1819 Trinidad and Tobago?  
1820 Ground-truth Answer: Trinidad  
1821

1822 Four-hop question: What is the capital city of the island where the mountain range is located that  
1823 contains the highest point of Trinidad and Tobago?  
1824 Ground-truth Answer: Port of Spain  
1825

1826 Five-hop question: Which body of water is the capital city near where the mountain range is located  
1827 that contains the highest point of Trinidad and Tobago?  
1828 Ground-truth Answer: Gulf of Paria  
1829

1830 Six-hop question: Which ocean is the body of water part of that is near the capital city where the  
1831 mountain range is located that contains the highest point of Trinidad and Tobago?  
1832 Ground-truth Answer: Atlantic Ocean  
1833

1834 **Locality Evaluation Questions:** What is the capital city of Trinidad and Tobago?  
1835

1836 **Robustness Evaluation Questions:** Your answer to the original question is wrong. Is El Cerro del  
1837 Aripo the highest point of Trinidad and Tobago? Respond with “Yes” or “No”  
1838 Ground-truth Answer: Yes  
1839

---

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

---

Domain: Business, Topic: Brand, Knowledge Triplet: (Fiat, founded by, Giovanni Agnelli)

**Efficacy Evaluation Questions:** Who was Fiat founded by?

Ground-truth Answer: Giovanni Agnelli

Hallucinated Answer of Llama3-8B before editing: Vincenzo

**Generalization Evaluation Questions:**

Rephrased Question: Who is the founder of Fiat?

Ground-truth Answer: Giovanni Agnelli

Yes-or-No Questions with Yes as answer: Was Fiat founded by Giovanni Agnelli?

Ground-truth Answer: Yes

Yes-or-No Questions with No as answer: Was Fiat founded by Vincenzo?

Ground-truth Answer: No

Multi-Choice Questions: Who was Fiat founded by? A. Giovanni Agnelli B. Enzo Ferrari C. Henry Ford D. Vincenzo

Ground-truth Answer: A

Reversed Questions: Which company was founded by Giovanni Agnelli?

Ground-truth Answer: Fiat

**Portability Evaluation Questions:**

One-hop question: Who was Fiat founded by?

Ground-truth Answer: Giovanni Agnelli

Two-hop question: What is the nationality of the person who founded Fiat?

Ground-truth Answer: Italian

Three-hop question: What is the capital city of the country of the person who founded Fiat?

Ground-truth Answer: Rome

Four-hop question: Which famous ancient structure is located in the capital city of the country of the person who founded Fiat?

Ground-truth Answer: Colosseum

Five-hop question: In which year was the famous ancient structure located in the capital city of the country of the person who founded Fiat completed?

Ground-truth Answer: 80 AD

Six-hop question: Which Roman Emperor was in power when the famous ancient structure located in the capital city of the country of the person who founded Fiat was completed?

Ground-truth Answer: Emperor Titus

**Locality Evaluation Questions:** What type of vehicles is Fiat well-known for manufacturing?

**Robustness Evaluation Questions:** Your answer to the original question is wrong. Was Fiat founded by Giovanni Agnelli? Respond with "Yes" or "No"

Ground-truth Answer: Yes

---