

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 different tasks. Meanwhile, knowledge editing has become a burgeoning paradigm that is designed to correct the erroneous factual knowledge encoded in LLMs for its advantage of avoiding retraining from scratch. However, one common issue of existing evaluation datasets for knowledge editing is that **they do not ensure 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?* Then, 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 assess 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 more future improvements and facilitate the progress in the field of knowledge editing. Data and code are available [here](#).

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cost of retraining LLMs from scratch, knowledge editing has been designed as a new paradigm to correct erroneous or outdated factual knowledge in LLMs (Wang et al., 2023).

Although there are many existing question-answering datasets such as WikiData_{recent} (Cohen et al., 2024), ZsRE (Yao et al., 2023), and WikiBio (Hartvigsen et al., 2024) widely used for knowledge editing evaluation, 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 being 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 of applying knowledge editing to LLMs. To better illustrate this point, following the evaluation setting in (Zhang et al., 2024), we conducted a preliminary study to examine the pre-edit and post-edit performances of Llama2-7B on the aforementioned three evaluation datasets. As shown in Table 1, we can clearly observe that Llama2-7B achieves a

Method	WikiData _{recent}	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.

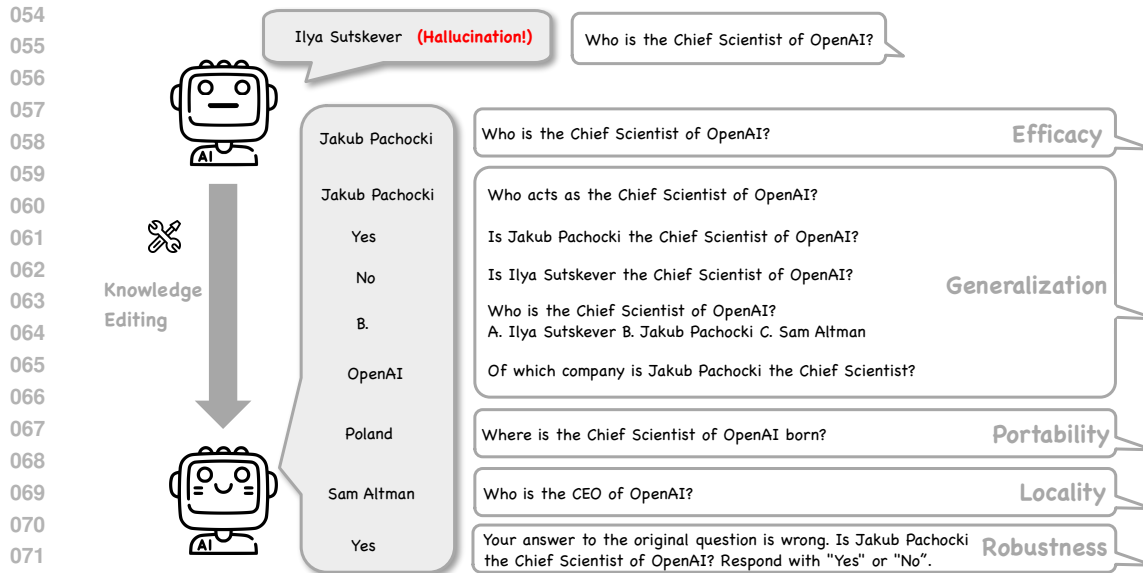


Figure 1: **Framework of HalluEditBench.** For real-world hallucinations, we holistically assess the performance of knowledge editing on *Efficacy*, *Generalization*, *Portability*, *Locality*, and *Robustness*.

relatively high performance, measured by the rate of answering the evaluation questions correctly (**Accuracy (%)**), even before applying knowledge editing techniques. Although the knowledge editing methods can bring an increase regarding 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 in correcting 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 Wikipedia. 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 (Meng et al., 2022), FT-M (Zhang et al., 2024), 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 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 marginally improve or negatively impact the *Generalization* performance. Editing techniques except ICE 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.

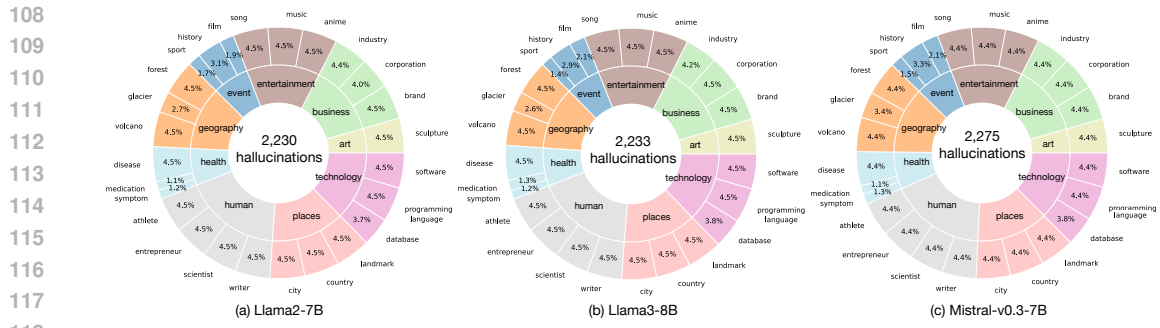


Figure 2: Statistics of HalluEditBench Across Topics and Domains.

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. An 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_{recent} (Cohen et al., 2024) and WikiBio (Hartvigsen et al., 2024)), we also choose Wikipedia as the factual knowledge source. In the *first* step, we retrieved 143,557 raw knowledge triplets using Wikidata Query Service (WDQS) from 26 topics, which can be categorized into 9 domains including *art*, *business*, *entertainment*, *event*, *geography*, *health*, *human*, *places*, and *technology*. 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 answers to questions about the object. When we construct factual questions and compare LLM-generated answers with the triplets, it would be hard 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., 2024c), 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 proposed to holistically assess the performance of knowledge editing methods in correcting hallucinations from five facets including *Efficacy*, *Gener-*

162 *alization, Portability, Locality, and Robustness*. First, we leveraged GPT-4o to generate evaluation
 163 question-answering pairs for each facet based on the hallucination dataset as well as the factuality
 164 verification questions in Section 2.1. One example of the generated evaluation QA pairs for each
 165 facet is shown in Figure 1. The specific prompt for GPT-4o is shown in Appendix C.

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

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

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

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

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

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

199 2.3 KNOWLEDGE EDITING TECHNIQUES

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

- 203 • **Locate-then-edit** is a popular knowledge editing paradigm that first locates factual knowledge at
 204 specific neurons or layers, and then makes modifications on them directly. We selected two typical
 205 methods ROME (Meng et al., 2022) and MEMIT (Meng et al., 2023) in HalluEditBench.
- 206 • **Fine-tuning** is a simple and straightforward way to update the parametric knowledge of LLMs. We
 207 selected three variations FT-L (Meng et al., 2022), FT-M (Zhang et al., 2024), and LoRA (Hu et al.,
 208 2022), which mitigate the catastrophic forgetting and overfitting issues of standard fine-tuning.
- 209 • **In-Context Editing** is a training-free paradigm that associates LLMs with in-context knowledge
 210 directly (Zheng et al., 2023; Shi et al., 2024; Fei et al., 2024). We adopted a simple baseline ICE
 211 method in (Zheng et al., 2023) that puts the new fact in the context and requires no demonstrations.
- 212 • **Memory-based** methods usually maintain a memory module for knowledge storage and updating.
 213 We selected a typical technique GRACE (Hartvigsen et al., 2024), which manages a discrete
 214 codebook and does not modify the original parameters. When encountering queries about edited
 215 knowledge, an adaptor adjusts layer-to-layer transformations with values searched in the codebook.

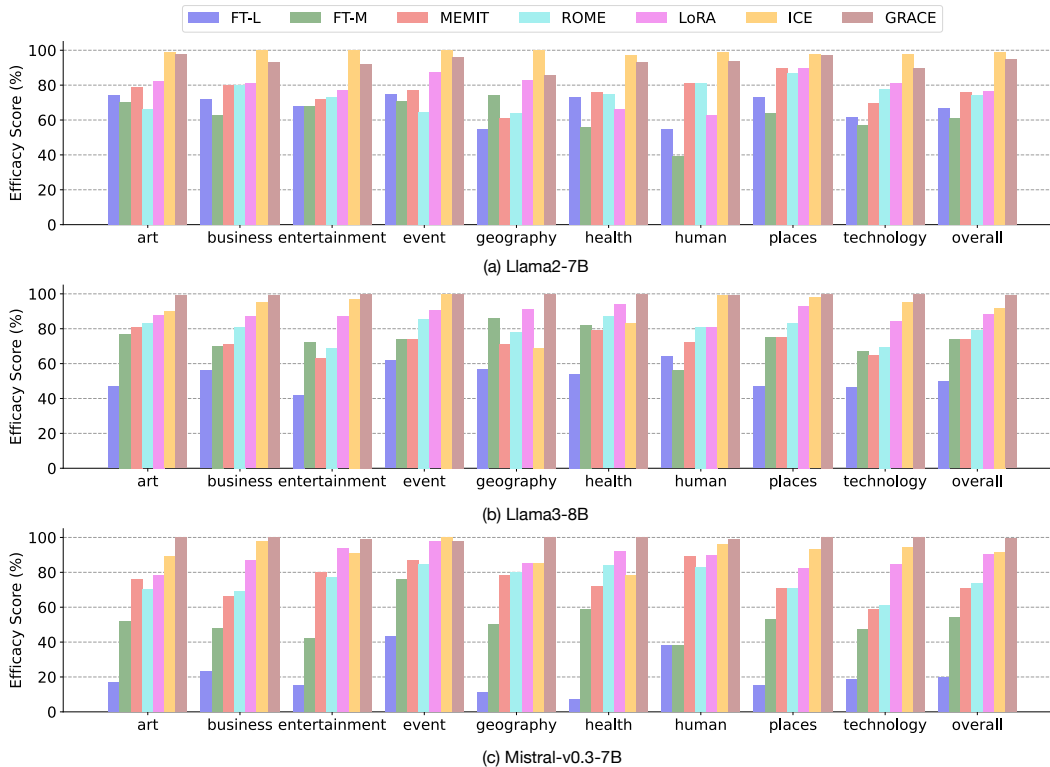


Figure 3: **Efficacy Scores of Knowledge Editing Methods.** The “overall” refers to the Efficacy Score (%) on the whole HalluEditBench embracing 9 domains for different methods. The Efficacy Score on each domain is also reported. Efficacy scores (%) are measured by the accuracy on Efficacy Evaluation 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* (in Appendix A).

3.1 FACET 1: EFFICACY

Figure 3 shows the Efficacy Score performance of post-edit LLMs in each domain and the whole HalluEditBench. Since we have ensured the LLMs’ pre-edit Efficacy Score is 0, Figure 3 can directly reflect the effectiveness of different knowledge editing techniques in correcting real-world hallucinations. Thus, we find that **the effectiveness of some techniques can be far from what their performance on previous datasets suggests**, implying the potential unreliability of performance on previous datasets. For example, as shown in Table 1, although FT-M achieves near 100% performance in existing datasets such as WikiData_{recent}, 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 original weights in LLMs, outperform the other methods, implying **the potential disadvantage of directly modifying parameters for editing knowledge**.

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. 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*.

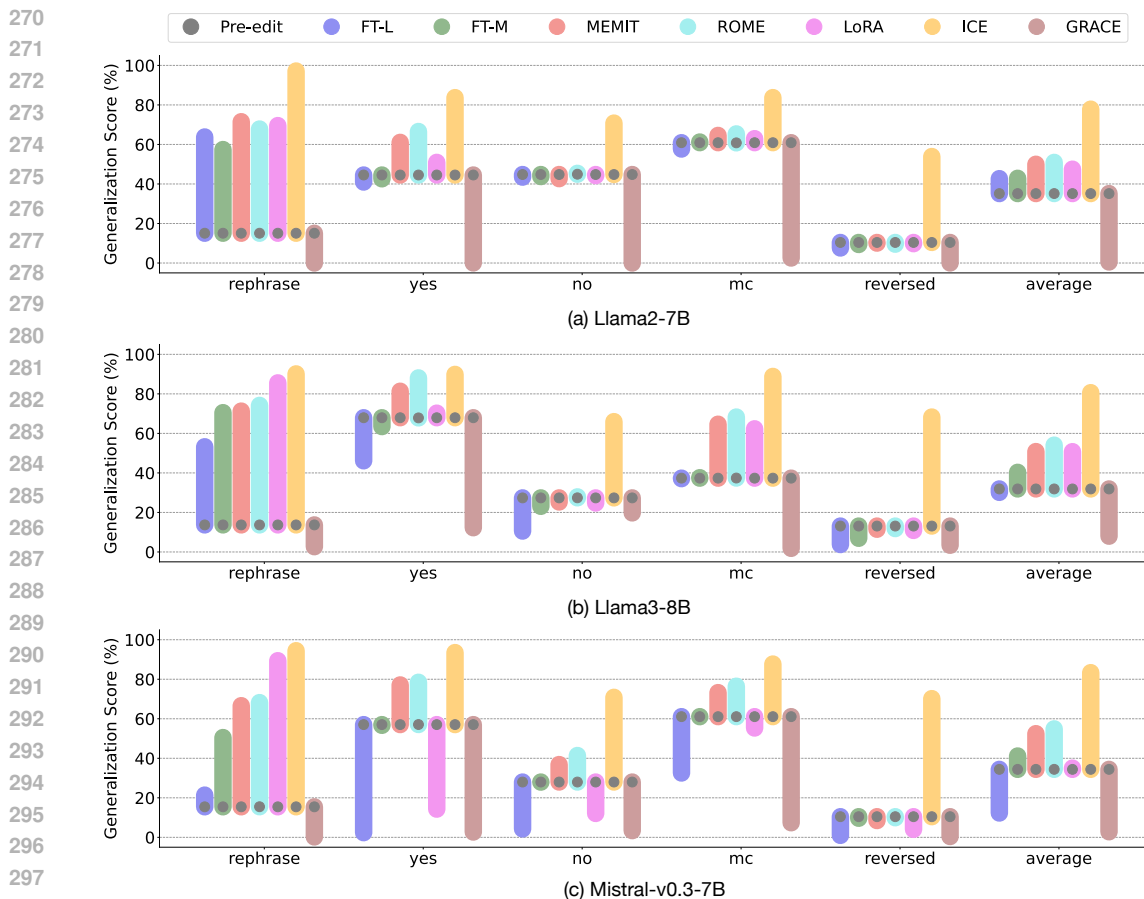


Figure 4: **Generalization Scores of Knowledge Editing Methods.** Generalization Scores (%) are measured by accuracy on five types of Generalization Evaluation Questions including Rephrased Questions (“rephrase”), Yes-or-No Questions with Yes or No as answers (“yes” or “no”), Multi-Choice Questions (“mc”), 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

Figure 4 shows both the pre-edit and post-edit Generalization Scores for different knowledge editing techniques on three LLMs. It is worth noting that the pre-edit performance is the same for different techniques and not zero, 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 marginally 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 marginally improve or negatively impact the *Generalization* performance.

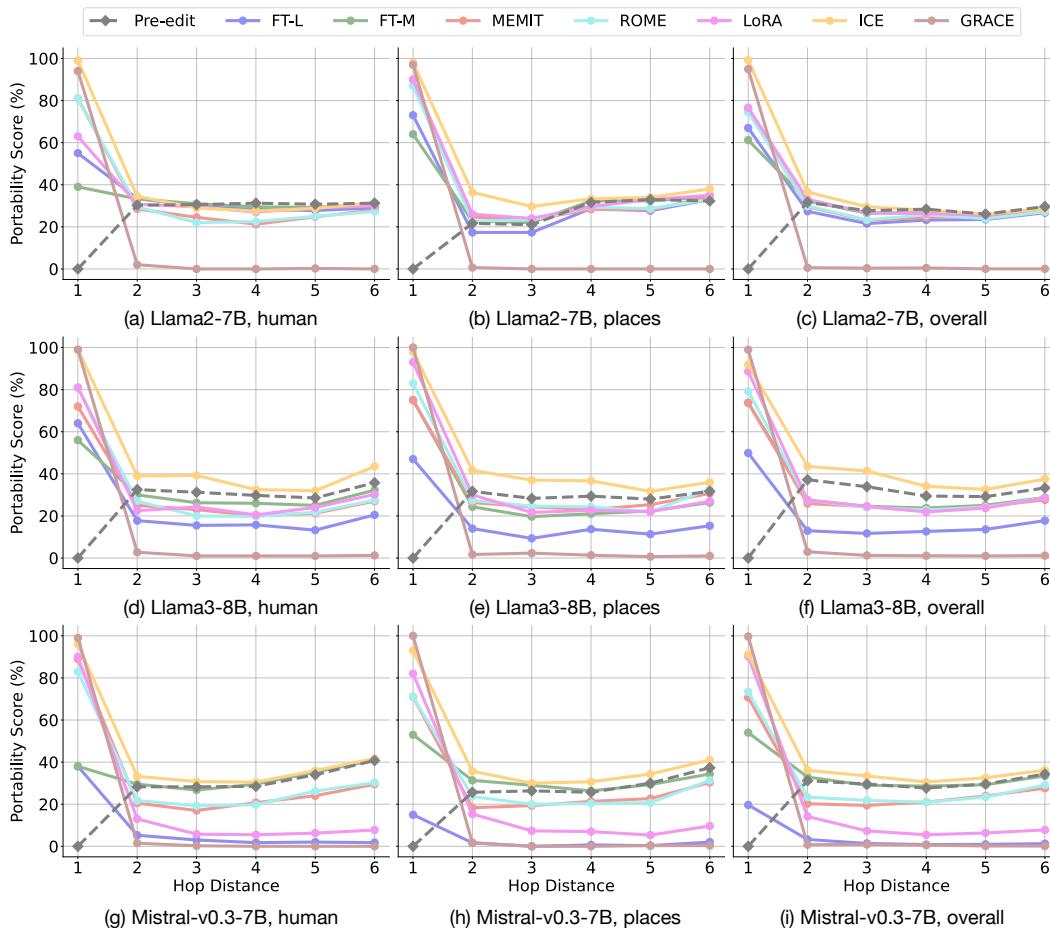


Figure 5: **Portability Scores of Knowledge Editing Methods.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions with N hops ($N = 1 \sim 6$). The Portability Evaluation Questions are the same as Efficacy Evaluation Questions 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 cannot really reason with the edited knowledge in multi-hop questions** regardless of 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 except ICE mostly underperform pre-edit Portability Scores, implying LLMs cannot really reason with edited knowledge in multi-hop questions.

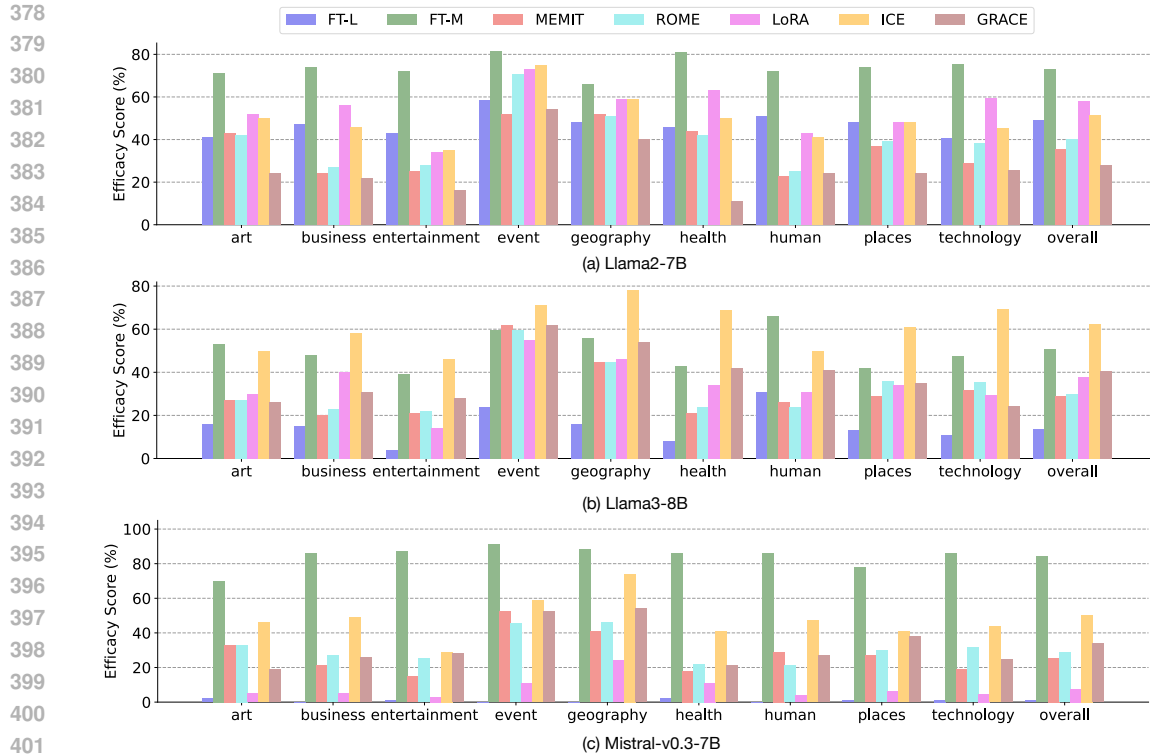


Figure 6: **Locality Scores of Knowledge Editing Methods.** Locality Scores (%) are measured by the unchanging rate on Locality Evaluation Questions after applying knowledge editing methods on LLMs. A higher Locality Score indicates that there is a higher percentage of LLMs’ answers to the unrelated 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*.

REFERENCES

- 432
433
434 Afra Feyza Akyürek, Eric Pan, Garry Kuwanto, and Derry Wijaya. Dune: Dataset for unified editing.
435 *ArXiv preprint*, abs/2311.16087, 2023. URL <https://arxiv.org/abs/2311.16087>.
- 436
437 Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects
438 of knowledge editing in language models. *Transactions of the Association for Computational*
439 *Linguistics*, 12:283–298, 2024.
- 440
441 Weizhi Fei, Xueyan Niu, Guoqing Xie, Yanhua Zhang, Bo Bai, Lei Deng, and Wei Han. Re-
442 trieval meets reasoning: Dynamic in-context editing for long-text understanding. *ArXiv preprint*,
443 abs/2406.12331, 2024. URL <https://arxiv.org/abs/2406.12331>.
- 444
445 Govind Gangadhar and Karl Stratos. Model editing by pure fine-tuning. *ArXiv preprint*,
446 abs/2402.11078, 2024. URL <https://arxiv.org/abs/2402.11078>.
- 447
448 Huaizhi Ge, Frank Rudzicz, and Zining Zhu. How well can knowledge edit methods edit perplexing
449 knowledge? *ArXiv preprint*, abs/2406.17253, 2024. URL <https://arxiv.org/abs/2406.17253>.
- 450
451 Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun
452 Peng. Model editing harms general abilities of large language models: Regularization to the rescue.
453 *ArXiv preprint*, abs/2401.04700, 2024. URL <https://arxiv.org/abs/2401.04700>.
- 454
455 Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi.
456 Aging with grace: Lifelong model editing with discrete key-value adaptors. *Advances in Neural*
457 *Information Processing Systems*, 36, 2024.
- 458
459 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
460 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth Interna-*
461 *tional Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.
462 OpenReview.net, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- 463
464 Han Huang, Haitian Zhong, Tao Yu, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. Vlkeb: A
465 large vision-language model knowledge editing benchmark. *arXiv preprint arXiv: 2403.07350*,
466 2024.
- 467
468 Jiaqi Li, Miaozeng Du, Chuanyi Zhang, Yongrui Chen, Nan Hu, Guilin Qi, Haiyun Jiang, Siyuan
469 Cheng, and Bozhong Tian. Mike: A new benchmark for fine-grained multimodal entity knowledge
470 editing. *ArXiv preprint*, abs/2402.14835, 2024a. URL <https://arxiv.org/abs/2402.14835>.
- 471
472 Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. Unveiling the
473 pitfalls of knowledge editing for large language models. In *The Twelfth International Conference*
474 *on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=fNktD3ib16>.
- 475
476 Zichao Li, Ines Arous, Siva Reddy, and Jackie Chi Kit Cheung. Evaluating dependencies in fact
477 editing for language models: Specificity and implication awareness. In *Findings of the Association*
478 *for Computational Linguistics: EMNLP 2023*, pp. 7623–7636, 2023.
- 479
480 Zihao Lin, Mohammad Beigi, Hongxuan Li, Yufan Zhou, Yuxiang Zhang, Qifan Wang, Wenpeng
481 Yin, and Lifu Huang. Navigating the dual facets: A comprehensive evaluation of sequential
482 memory editing in large language models. *ArXiv preprint*, abs/2402.11122, 2024. URL <https://arxiv.org/abs/2402.11122>.
- 483
484 Zeyu Leo Liu, Shrey Pandit, Xi Ye, Eunsol Choi, and Greg Durrett. Codeupdatearena: Benchmarking
485 knowledge editing on api updates. *ArXiv preprint*, abs/2407.06249, 2024. URL <https://arxiv.org/abs/2407.06249>.
- Jun-Yu Ma, Jia-Chen Gu, Zhen-Hua Ling, Quan Liu, and Cong Liu. Untying the reversal curse
via bidirectional language model editing. *ArXiv preprint*, abs/2310.10322, 2023. URL <https://arxiv.org/abs/2310.10322>.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.

- 486 Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing
487 memory in a transformer. In *The Eleventh International Conference on Learning Representations*,
488 2023. URL <https://openreview.net/forum?id=MkbcAHIYgyS>.
- 489
490 Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. Memory-
491 based model editing at scale. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári,
492 Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022,*
493 *17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning*
494 *Research*, pp. 15817–15831. PMLR, 2022. URL [https://proceedings.mlr.press/v162/
495 mitchell122a.html](https://proceedings.mlr.press/v162/mitchell122a.html).
- 496 Derek Powell, Walter Gerych, and Thomas Hartvigsen. Taxi: Evaluating categorical knowledge
497 editing for language models. *ArXiv preprint*, abs/2404.15004, 2024. URL [https://arxiv.org/
498 abs/2404.15004](https://arxiv.org/abs/2404.15004).
- 499 Domenic Rosati, Robie Gonzales, Jinkun Chen, Xuemin Yu, Melis Erkan, Yahya Kayani,
500 Satya Deepika Chavatapalli, Frank Rudzicz, and Hassan Sajjad. Long-form evaluation of model
501 editing. *ArXiv preprint*, abs/2402.09394, 2024. URL <https://arxiv.org/abs/2402.09394>.
- 502
503 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askeel, Samuel R. Bowman,
504 Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam
505 McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and
506 Ethan Perez. Towards understanding sycophancy in language models. In *The Twelfth International*
507 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?id=
508 tvhaxkMKAn](https://openreview.net/forum?id=tvhaxkMKAn).
- 509 Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu.
510 Retrieval-enhanced knowledge editing for multi-hop question answering in language models. *ArXiv*
511 *preprint*, abs/2403.19631, 2024. URL <https://arxiv.org/abs/2403.19631>.
- 512
513 Haoyu Wang, Tianci Liu, Tuo Zhao, and Jing Gao. Roselora: Row and column-wise sparse low-rank
514 adaptation of pre-trained language model for knowledge editing and fine-tuning. *ArXiv preprint*,
515 abs/2406.10777, 2024a. URL <https://arxiv.org/abs/2406.10777>.
- 516 Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang,
517 and Huajun Chen. Wise: Rethinking the knowledge memory for lifelong model editing of large
518 language models. *ArXiv preprint*, abs/2405.14768, 2024b. URL [https://arxiv.org/abs/2405.
519 14768](https://arxiv.org/abs/2405.14768).
- 520
521 Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, et al. Knowledge editing
522 for large language models: A survey. *ArXiv preprint*, abs/2310.16218, 2023. URL <https://arxiv.org/abs/2310.16218>.
- 523
524 Wenxuan Wang, Juluan Shi, Zhaopeng Tu, Youliang Yuan, Jen-tse Huang, Wenxiang Jiao, and
525 Michael R Lyu. The earth is flat? unveiling factual errors in large language models. *ArXiv preprint*,
526 abs/2401.00761, 2024c. URL <https://arxiv.org/abs/2401.00761>.
- 527
528 Yifan Wei, Xiaoyan Yu, Huanhuan Ma, Fangyu Lei, Yixuan Weng, Ran Song, and Kang Liu. Assess-
529 ing knowledge editing in language models via relation perspective. *ArXiv preprint*, abs/2311.09053,
530 2023. URL <https://arxiv.org/abs/2311.09053>.
- 531
532 Zihao Wei, Jingcheng Deng, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. Mlake: Mul-
533 tilingual knowledge editing benchmark for large language models. *ArXiv preprint*, abs/2404.04990,
534 2024. URL <https://arxiv.org/abs/2404.04990>.
- 535
536 Suhang Wu, Minlong Peng, Yue Chen, Jinsong Su, and Mingming Sun. Eva-kellm: A new benchmark
537 for evaluating knowledge editing of llms. *ArXiv preprint*, abs/2308.09954, 2023. URL <https://arxiv.org/abs/2308.09954>.
- 538
539 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen,
and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. *ArXiv*
preprint, abs/2305.13172, 2023. URL <https://arxiv.org/abs/2305.13172>.

- 540 Lang Yu, Qin Chen, Jie Zhou, and Liang He. Melo: Enhancing model editing with neuron-indexed
541 dynamic lora. *ArXiv preprint*, abs/2312.11795, 2023. URL [https://arxiv.org/abs/2312.](https://arxiv.org/abs/2312.11795)
542 [11795](https://arxiv.org/abs/2312.11795).
- 543
544 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi,
545 Shengyu Mao, Jintian Zhang, Yuansheng Ni, et al. A comprehensive study of knowledge editing
546 for large language models. *ArXiv preprint*, abs/2401.01286, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2401.01286)
547 [abs/2401.01286](https://arxiv.org/abs/2401.01286).
- 548 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lema Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,
549 Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi.
550 Siren’s song in the ai ocean: A survey on hallucination in large language models. *ArXiv preprint*,
551 abs/2309.01219, 2023. URL <https://arxiv.org/abs/2309.01219>.
- 552
553 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
554 Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen,
555 Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and
556 Ji-Rong Wen. A survey of large language models. *ArXiv preprint*, abs/2303.18223, 2023. URL
557 <https://arxiv.org/abs/2303.18223>.
- 558 Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can
559 we edit factual knowledge by in-context learning? In Houda Bouamor, Juan Pino, and Kalika Bali
560 (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*,
561 pp. 4862–4876, Singapore, 2023. Association for Computational Linguistics. doi: 10.18653/v1/
562 2023.emnlp-main.296. URL <https://aclanthology.org/2023.emnlp-main.296>.
- 563 Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen.
564 Mquake: Assessing knowledge editing in language models via multi-hop questions. *ArXiv*
565 *preprint*, abs/2305.14795, 2023. URL <https://arxiv.org/abs/2305.14795>.
- 566
567 Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv
568 Kumar. Modifying memories in transformer models. *ArXiv preprint*, abs/2012.00363, 2020. URL
569 <https://arxiv.org/abs/2012.00363>.
- 570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Content of Appendix

A Facet 5: Robustness	13
B Related Work	14
C Reproducibility Statement	15
D Details of the Benchmarked Knowledge Editing Techniques	16
E More Experiment Results	18
E.1 Generalization Scores of Knowledge Editing Methods on Each Domain	18
E.2 Portability Scores of Knowledge Editing Methods on More Domains	23
E.3 Robustness Scores of Knowledge Editing Methods on More Domains	26

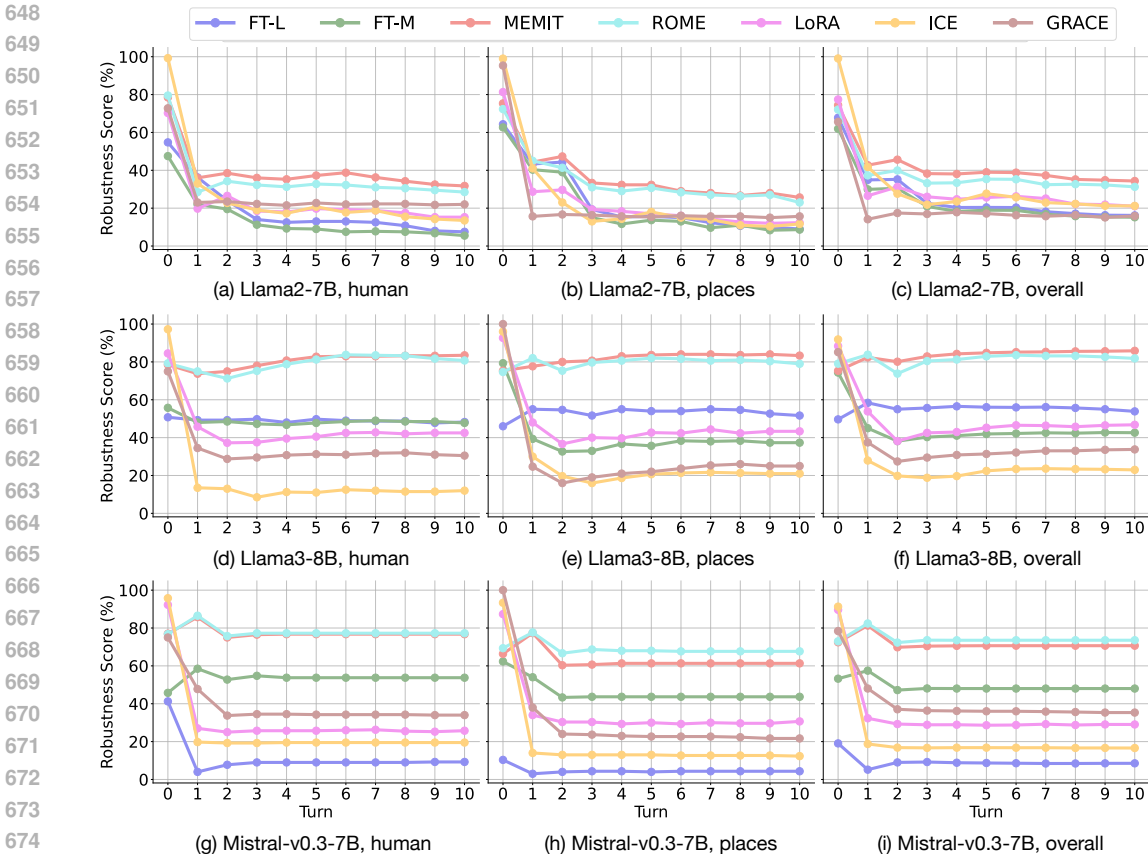


Figure 7: **Robustness Scores of Knowledge Editing Methods.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with M turns ($M = 1 \sim 10$). We regard Efficacy Scores as the Robustness Scores when M is 0. The Robustness Scores on two domains “human” and “places” are 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.

A 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 technique 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, FT-M on Llama3-8B and MEMIT, ROME, FT-M, FT-L on Mistral-v0.3-7B maintain almost the same performance as turns increase, suggesting a relatively 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*.

B RELATED WORK

Knowledge editing techniques have attracted increasing attention for their efficiency advantages in addressing obsolete or hallucinated information in LLMs (Wang et al., 2023; Zhang et al., 2024). 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., 2024b; 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; Wei et al., 2023; 2024; Ge et al., 2024; Huang et al., 2024; Liu et al., 2024; Ma et al., 2023; Li et al., 2024a;b; 2023; Zhong et al., 2023; Wu et al., 2023; Powell et al., 2024; Lin et al., 2024; Akyürek et al., 2023; Gu 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. (2024) 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.

C REPRODUCIBILITY STATEMENT

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

- Llama3-8b: <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 following prompt to generate evaluation questions in *Generalization* and *Locality* aspects.

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 and a relation, 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 avoid including any correct answers from other multi-hop questions.

Ensure the answers for multi-hop questions are correct, and do not use ‘N/A’ as answers. Output in JSON format. Below is an example:

Example input:

subject: Amazon, relation: founder

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",
}
```

D DETAILS OF THE BENCHMARKED KNOWLEDGE EDITING TECHNIQUES

FT-L (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., 2024) 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.

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.

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.

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.

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.

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

875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

E MORE EXPERIMENT RESULTS

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

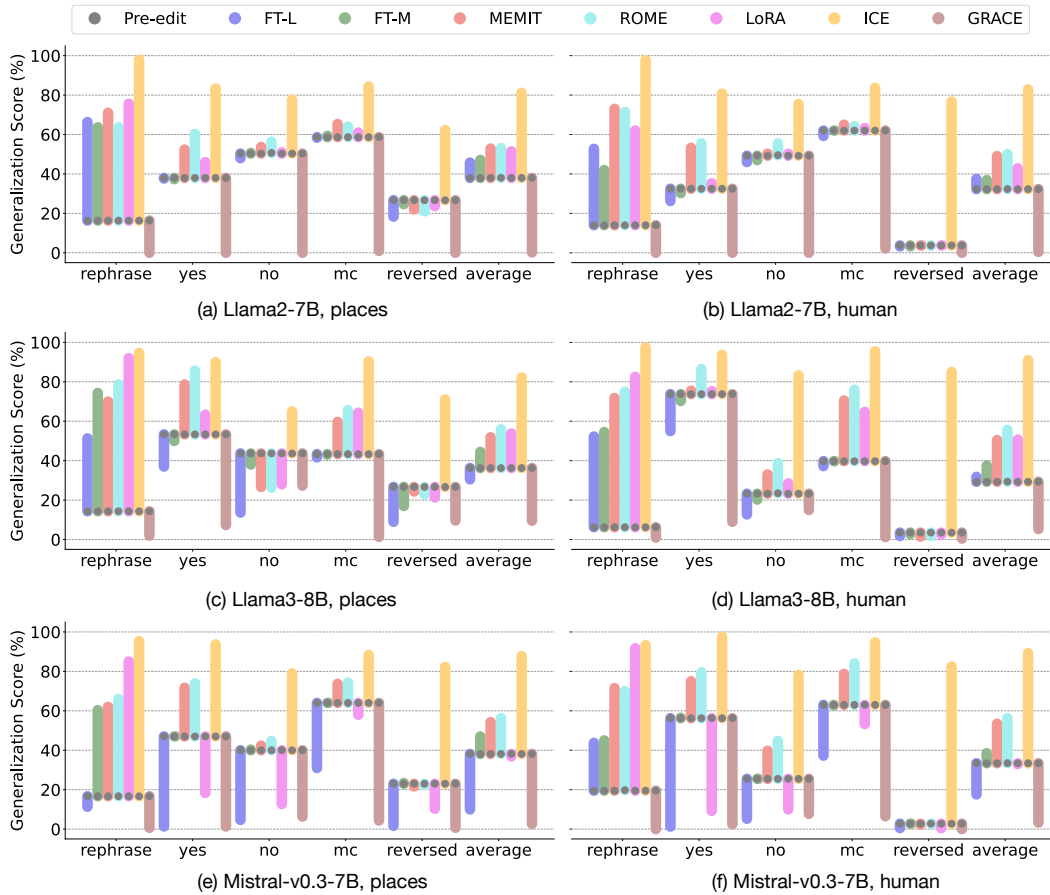


Figure 8: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No 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 “places” and “human”.

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

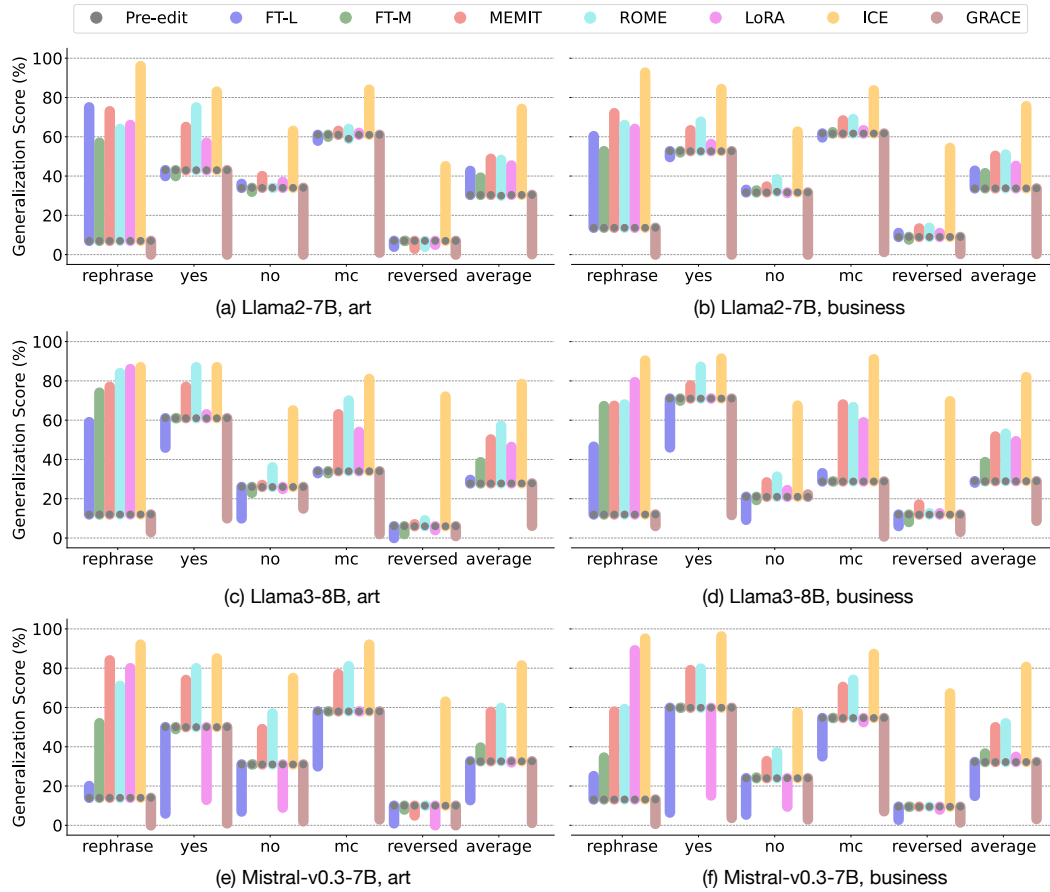


Figure 9: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No 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”.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

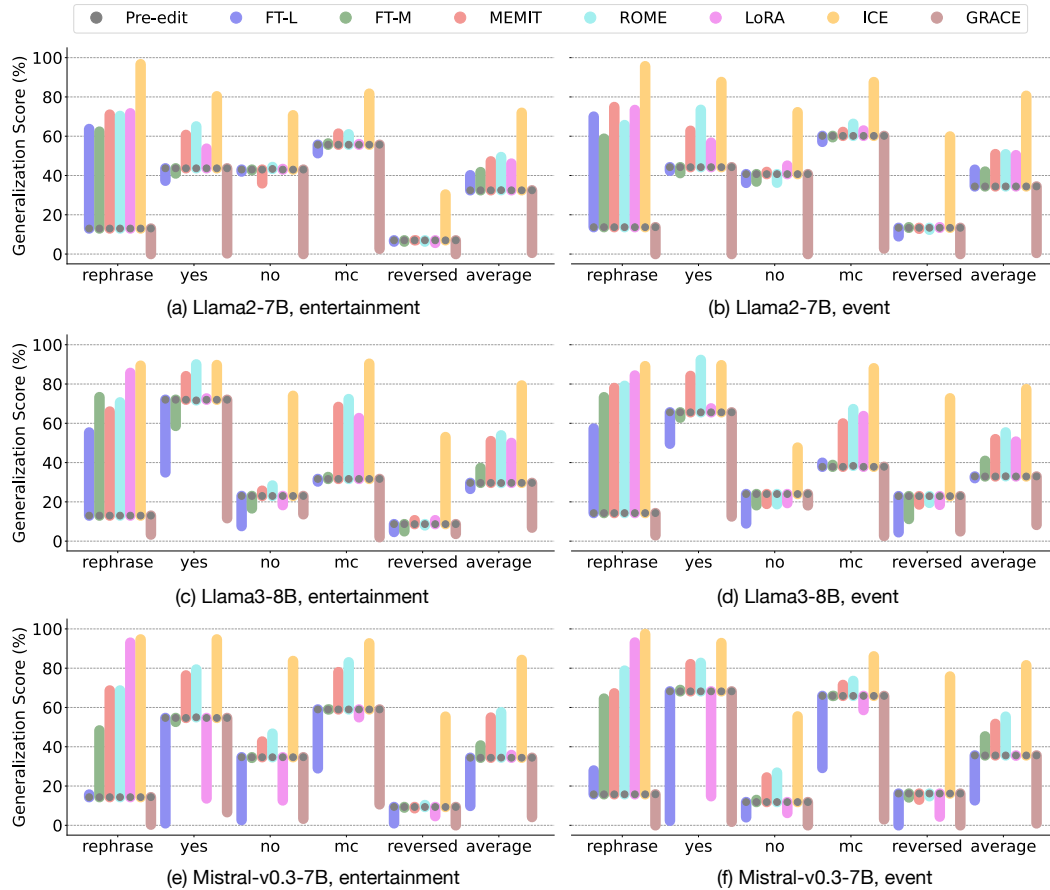


Figure 10: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No 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”.

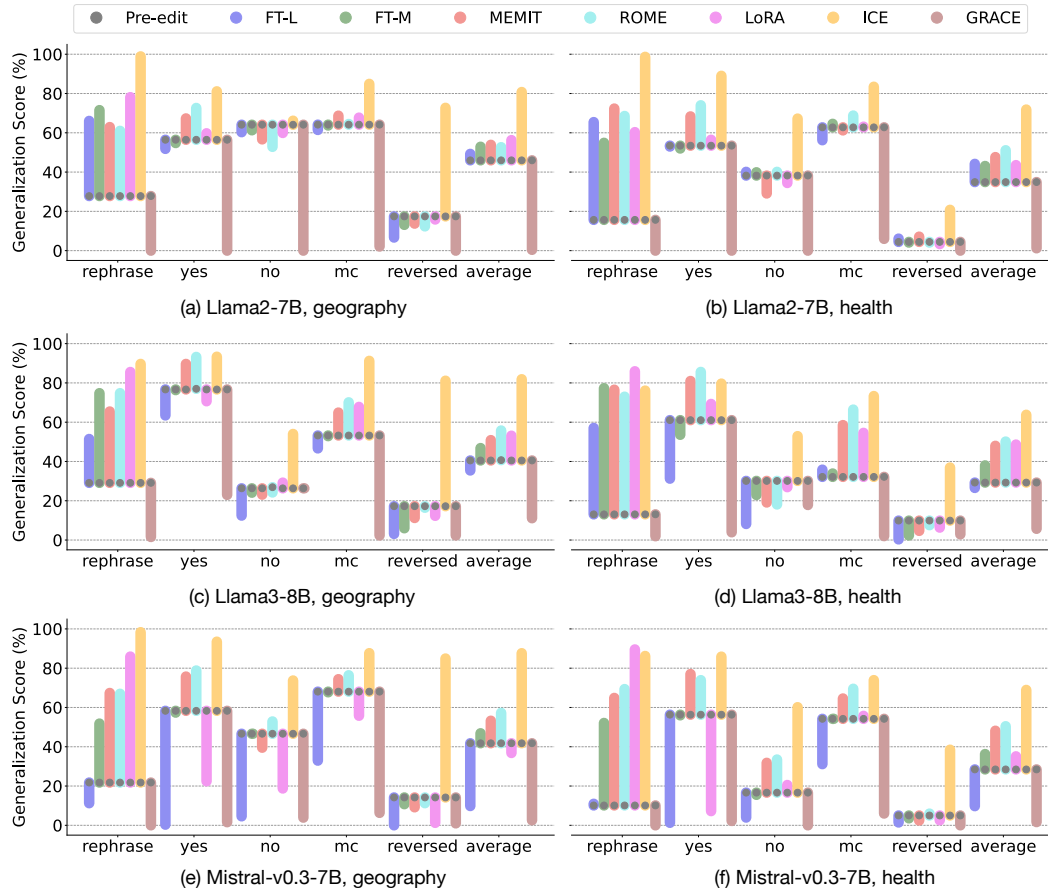


Figure 11: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No 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”.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

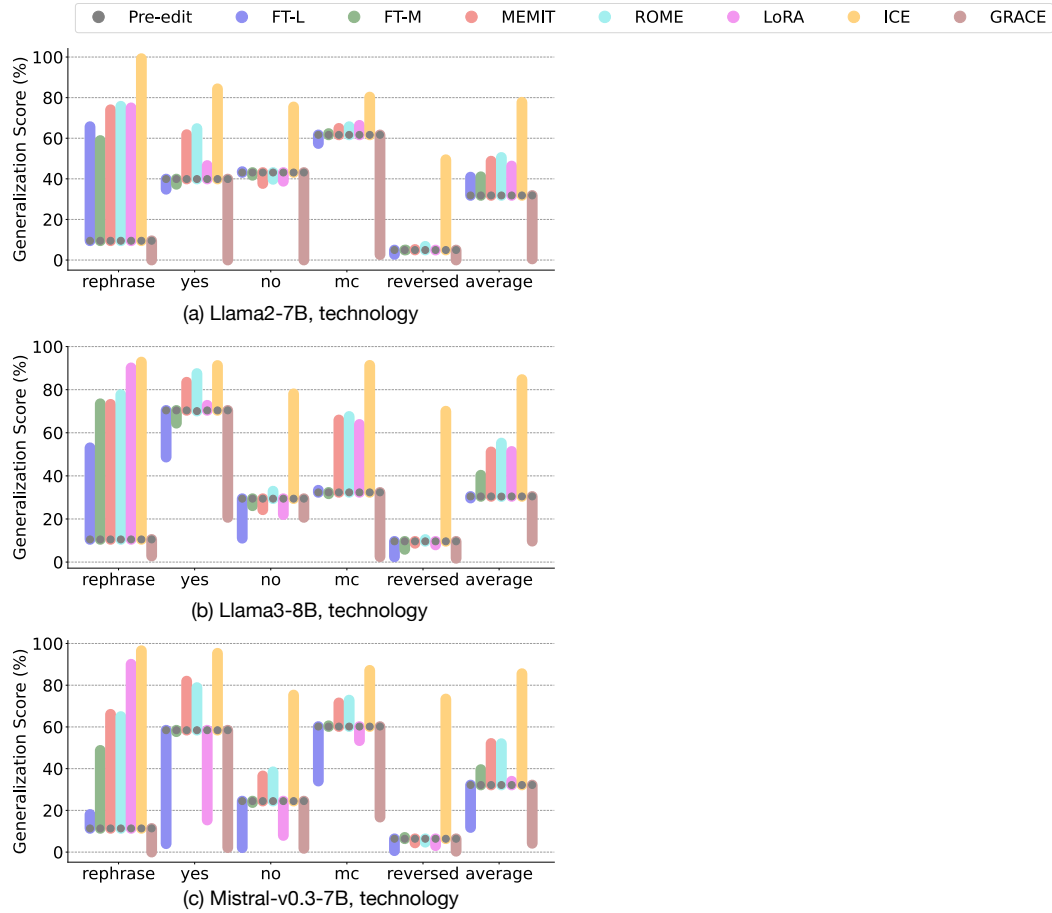


Figure 12: **Generalization Scores of Knowledge Editing Methods on 3 LLMs and 2 Domains.** Generalization Scores (%) are measured by the accuracy on five types of Generalization Evaluation Question-answer Pairs including Rephrased Questions (“rephrase”), two types of Yes-or-No Questions with Yes or No 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

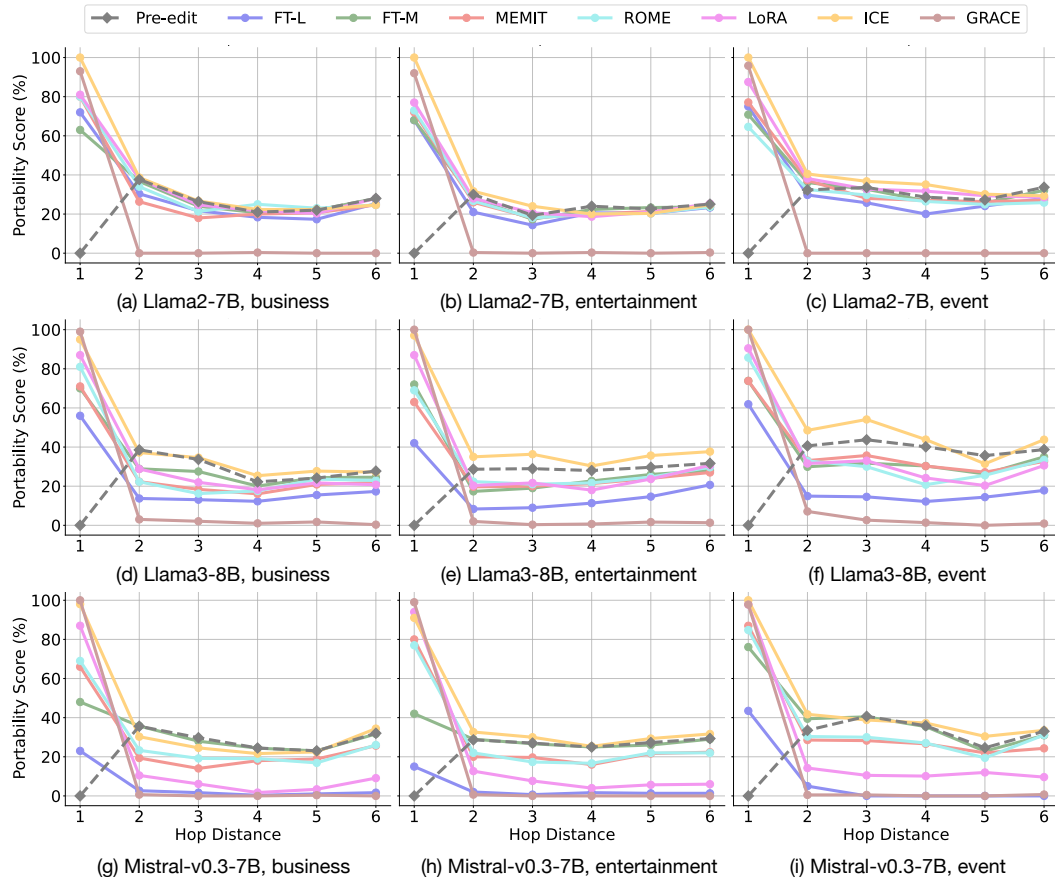


Figure 13: **Portability Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions when with N hops. The Portability Evaluation Questions are the same as Efficacy Evaluation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

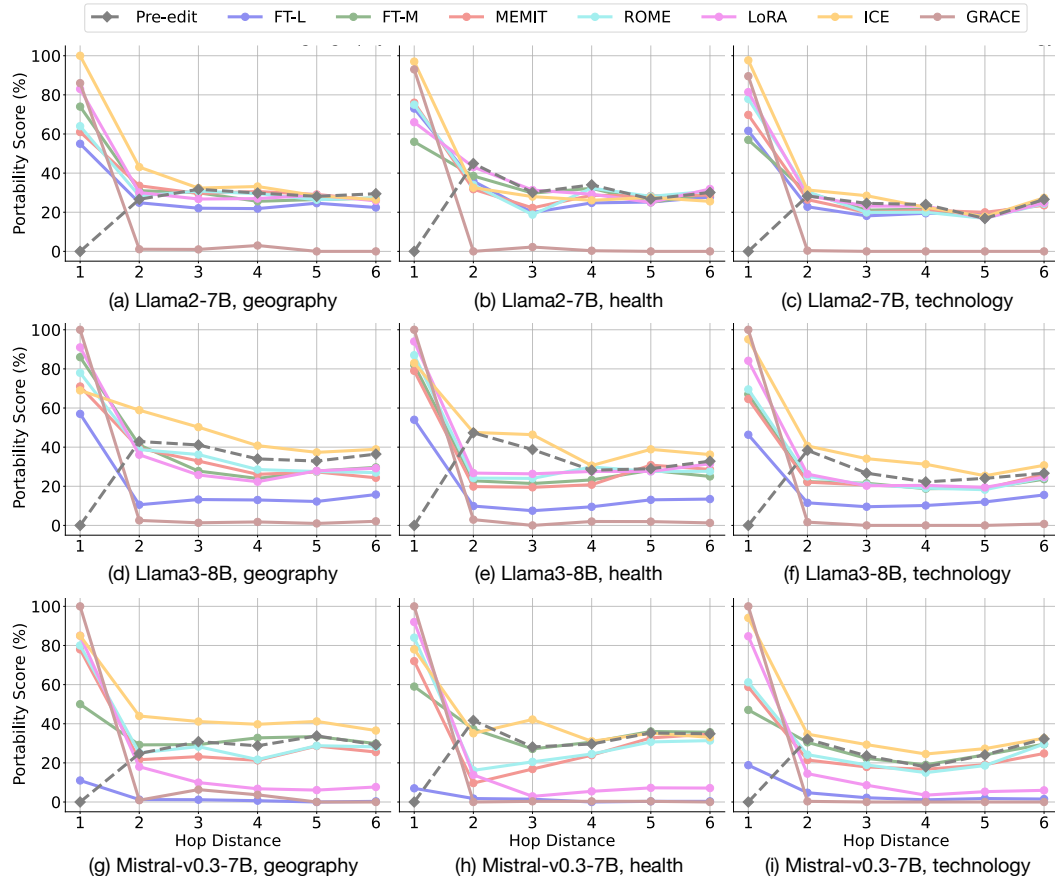


Figure 14: **Portability Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions when with N hops. The Portability Evaluation Questions are the same as Efficacy Evaluation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

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

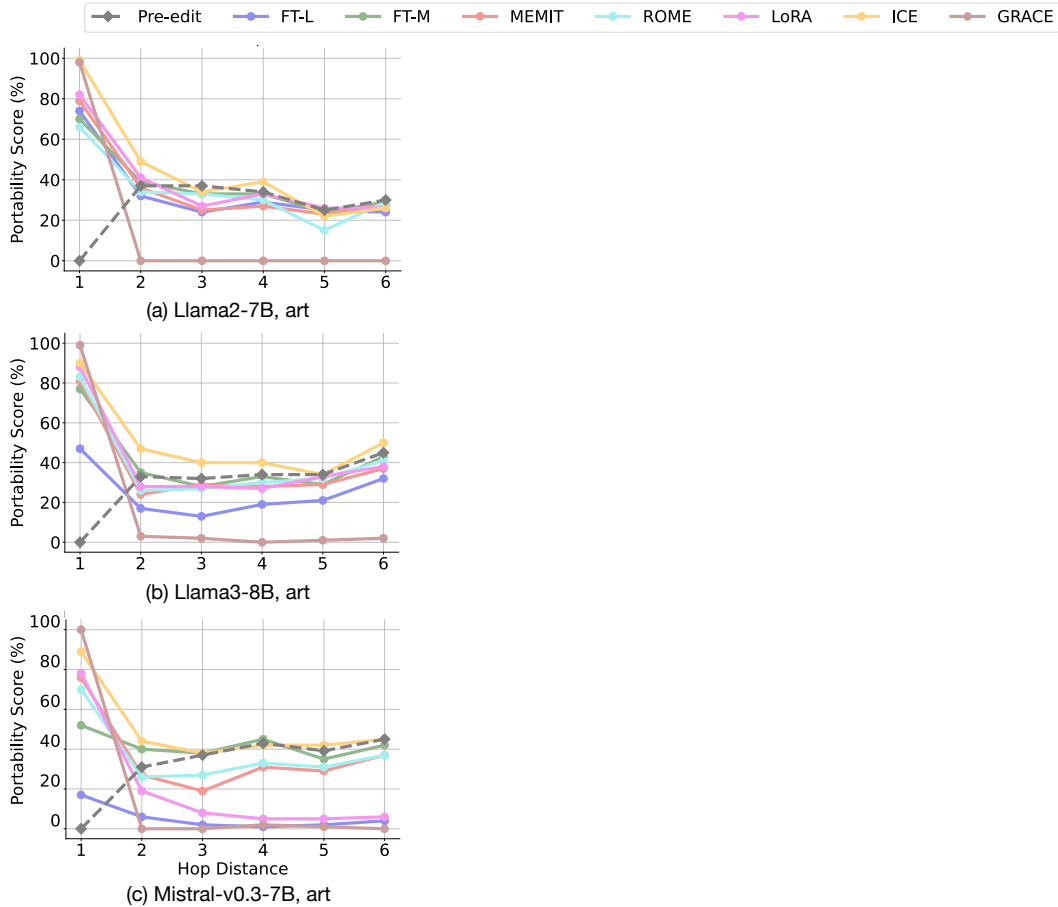


Figure 15: **Portability Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Portability Scores (%) are measured by the accuracy on Portability Evaluation Questions, which are Efficacy Evaluation Questions when with N hops. The Portability Evaluation Questions are the same as Efficacy Evaluation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

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

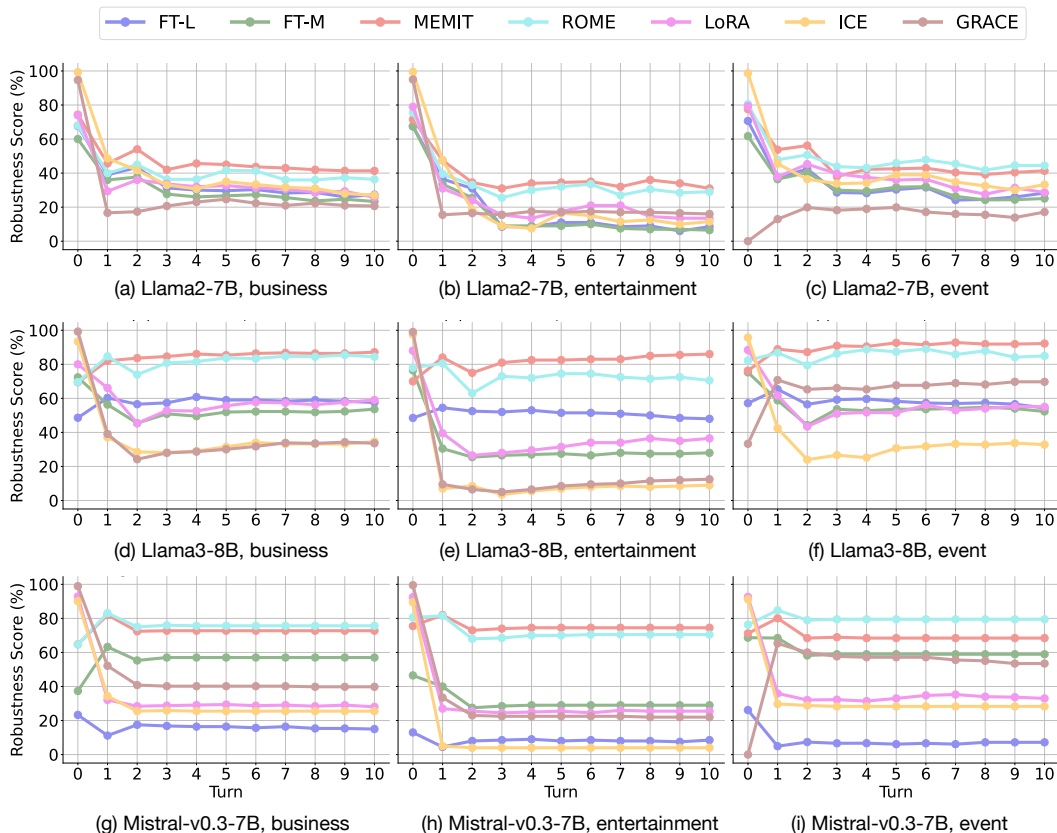


Figure 16: **Robustness Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with M turns ($M = 1 \sim 10$). We regard Efficacy Scores as the Robustness Scores when M is 0. The domains include “business”, “entertainment”,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

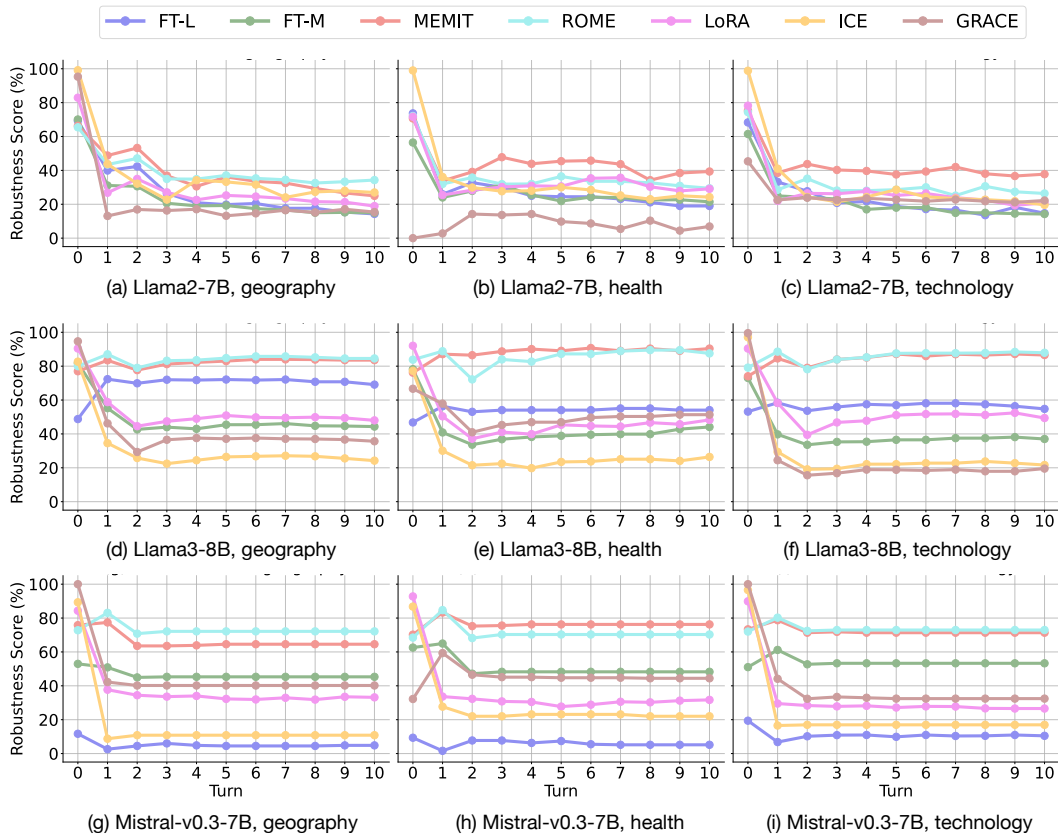


Figure 17: **Robustness Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with M turns ($M = 1 \sim 10$). We regard Efficacy Scores as the Robustness Scores when M is 0. The domains include “geography”, “health”,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

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

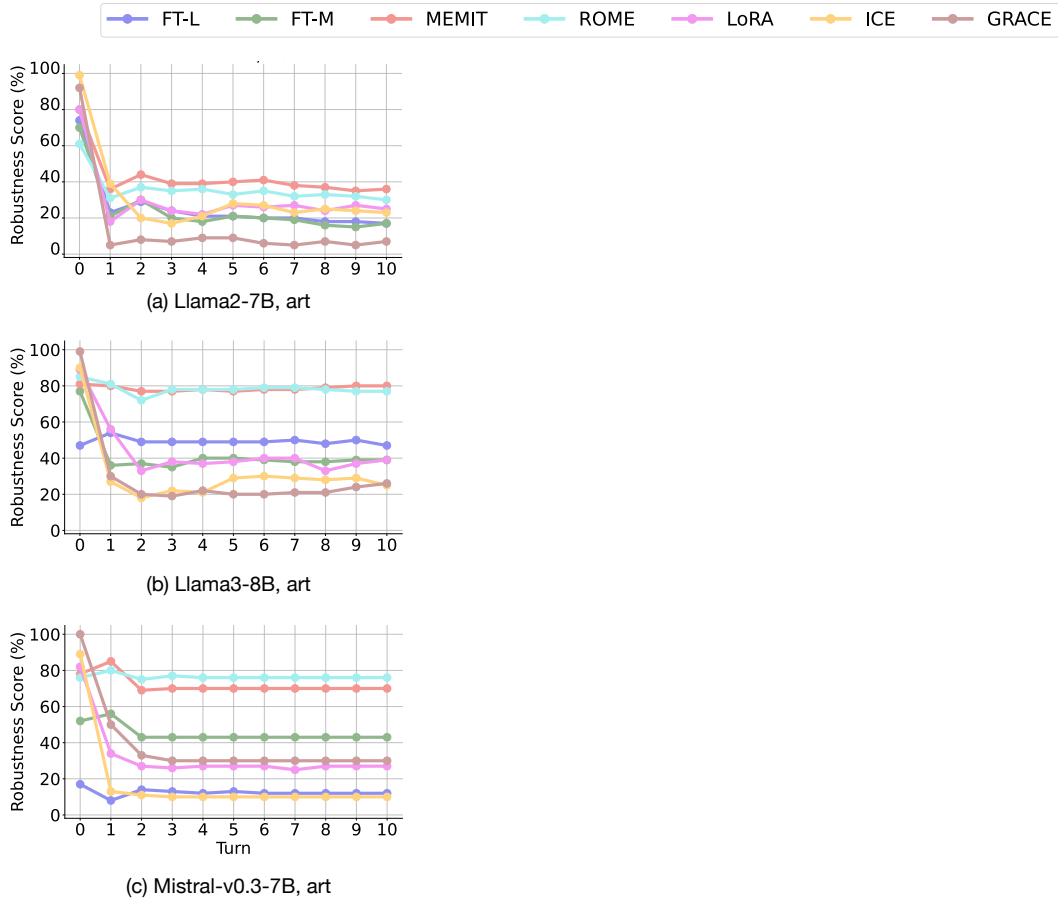


Figure 18: **Robustness Scores of Knowledge Editing Methods on 3 LLMs and 3 Domains.** Robustness Scores are calculated by the accuracy on Robustness Evaluation Questions with M turns ($M = 1 \sim 10$). We regard Efficacy Scores as the Robustness Scores when M is 0. The domain is “art”.