

Rebuilding ROME : Resolving Model Collapse during Sequential Model Editing

Akshat Gupta¹, Sidharth Baskaran², Gopala Anumanchipalli¹

¹UC Berkeley, ²Automorphic Inc.

akshat.gupta@berkeley.edu, sid@automorphic.ai

Abstract

Recent work using Rank-One Model Editing (ROME), a popular model editing method, has shown that there are certain facts that the algorithm is unable to edit without breaking the model. Such edits have previously been called disabling edits (Gupta et al., 2024a). These disabling edits cause immediate model collapse and limits the use of ROME for sequential editing. In this paper, we show that disabling edits are an artifact of irregularities in the implementation of ROME. With this paper, we provide a more stable implementation ROME, which we call r-ROME and show that model collapse is no longer observed when making large scale sequential edits with r-ROME, while further improving generalization and locality of model editing compared to the original implementation of ROME.

1 Introduction

Large language models (LLMs) are expensive to train and the knowledge contained in these models gets obsolete with time. Model editing or knowledge editing (Yao et al., 2023) has recently come out as a popular method to update knowledge in large language models (LLMs). In this paper, we focus on one popular parameter-modifying model editing methods called ROME (Rank-One Model Editing) (Meng et al., 2022a). ROME is not only one of the most popular model editing algorithms, but is also widely used in unlearning (Patil et al., 2023) and model interpretability (Ghandeharioun et al., 2024; Geva et al., 2023) literature.

While a lot of model editing approaches perform well when making singular edits, editing multiple facts in a model still remains a challenge for parameter-modifying model editing methods. One way to make multiple edits to the same model is through **sequential editing** (Yao et al., 2023) - where we make a series of single edits to a model by modifying the parameters of the model after

After DISABLING EDIT:

One can get to Bay by navigating the, B Serie Italian, as. Italian Serie Gi Italy Italy Italy Italy Italy Italy Italy is Italian,\n Italy mag Tur_____

After Normal EDIT:

The language used by people in Equatorial Guinea is often very informal and it is important to be aware of what is acceptable and not acceptable. This is particularly important if there has been a _____

Figure 1: A typical generation example after a disabling edit is compared to a normal model edit using ROME. The bold and underlined part in the text is input prompt.

every edit. Recent works have started studying the effects of sequential editing and found that ROME (Meng et al., 2022a) was prone to a sudden model collapse by a single edit (Gupta et al., 2024a; Yang et al., 2024; Hu et al., 2024). This effect was first observed in Gupta et al. (2024a) during sequential editing. The collapse included complete loss of downstream performance, inability to recall previously editing facts and loss of the ability to even get edited. Such facts were named **disabling edits** by Gupta et al. (2024a) and were later independently observed by Yang et al. (2024); Hu et al. (2024).

Disabling edits are detrimental for knowledge editing at scale. While a gradual model degradation is expected as we make sequential edits to a model (Gupta et al., 2024a), disabling edits lead to a sudden model collapse irrespective of when the disabling fact is edited, making sequential editing impossible. An example of this can be seen in Figure 3a, where instead of allowing gradual model degradation when doing sequential editing like in Figure 4, the presence of disabling edits lead to a sudden and immediate model collapse.

In this paper, we aim to find the source of these disabling edits. We first introduce two metrics for identifying disabling edits - generation entropy and the norm of matrix update. We plot edits made by ROME along these two dimensions and show new ways of identifying disabling edits even when

DATASET	IMPLEMENTATION	Efficacy		Generalization		Locality		Score
		ES \uparrow	EM \uparrow	PS \uparrow	PM \uparrow	NS \uparrow	NM \uparrow	S \uparrow
CF	ORIGINAL	99.92	99.68	96.29	71.58	75.8	10.25	89.32
	r-ROME	99.74	97.79	99.09	70.86	80.62	26.0	92.22
	p-ROME	99.9	99.36	97.04	63.01	80.0	5.74	91.42

Table 1: The above represents model editing results for 5000 singular model edits made on GPT-J-6B from the CounterFact dataset (non-sequential).

making singular edits. As we dig deeper into the optimization objectives and the codebase of ROME, we find that the disabling edits in ROME are a result of irregularities in the implementation of ROME, and not an artifact of the optimization objective. Specifically, disabling edits were caused due to the asymmetric usage of key-vectors in the update equation of ROME. With this paper, we share our new ROME code-base and invite researchers to use it for model editing. Our implementation of ROME, which we call r-ROME, can be found [here](#)¹.

2 Background

Facts are usually added in ROME using key-value format, where a key is the vector representation of a query-phrase and the value is the vector representation of the target object. For example, when adding a new fact - "*The president of USA is John Cena*", the query-phrase here is "*The president of USA is*" and the target object is "*John Cena*". The key-vector is defined by [Meng et al. \(2022a\)](#) is the activation of the first linear layer in the MLP targeted by ROME:

$$k^{(l^*)}(x) = \sigma \left(W_{fc}^{(l^*)} \gamma \left(a_{[x],i}^{(l^*)} + h_{[x],i}^{(l^*-1)} \right) + b_{fc}^{(l^*)} \right) \quad (1)$$

Editing in ROME is done using a pair of vectors - (k_e, v_e) that represent a new fact being added. k_e , also called the key-vector is a vector representation of the query-phrase, and v_e , or the value-vector is the vector representation of the target object. The weights of the specific layer being edited in ROME are updated from W_0 to \hat{W} by inserting a new fact (k_e, v_e) using the following equation:

$$\hat{W} = W_0 + \Delta \quad (2)$$

where $\Delta = (v_e - W_0 k_e) \frac{k_e^T C_0^{-1}}{k_e^T C_0^{-1} k_e}$

where Δ is the update to the current weight matrix being edited such that the new fact (k_e, v_e) gets incorporated. Additionally, each key-vector in k_e is not just the representation of a single prompt. To enhance generalization, [Meng et al. \(2022a,b\)](#) create the key-vector as an average representations over the query-phrase with random prefixes. This is done so that the represented key-vectors do not just represent one way to phrase the query-phrase and edits made using these representations can generalize over different paraphrases of the edited facts. The final key vector is found by averaging over N random prefixes using the equation:

$$k_e = \frac{1}{N} \sum_{i=1}^N k(x_i \oplus p) \quad (3)$$

Here $k(x_i \oplus p)$ represents the key-vector corresponding to a prefix x_i being concatenated with the original query-phrase p . Examples of prefixes added in ROME can be seen in Table 3. In this paper, we will refer to the averaged prefix representation of keys with k_e , whereas when the representation just consists of the original prompt, we will depict that with a superscript as k_e^o . The following equation explicitly differentiates between the two mathematically:

$$k_e^o = k(p) \quad (4)$$

Evaluating Model Editing. Model editing is usually evaluated along three metrics - reliability, generalization and locality. Reliability represents if a fact was successfully added in a model and is measured using edit score (ES) and edit magnitude (EM) metrics. ES measures the portion of cases when an edited fact is more probable than the original fact post-editing, whereas EM measures the difference in the probability magnitudes of the edited and original facts. Generalization represents if the edited fact is recalled through paraphrases of the prompt used to edit the fact and is measured

¹<https://github.com/scalable-model-editing/rebuilding-rome>

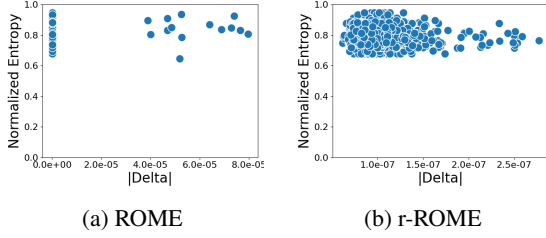


Figure 2: This figure shows the difference between the ROME and r-ROME updates on GPTJ (6B) for 5k individual edits. Our implementation shows much less potential disabling edits indicated by lower $|\Delta|$ values.

using paraphrase score (PS) and paraphrase magnitude defined similarly as above for paraphrases of the edited facts. Locality represents if editing of one fact affects other facts stored inside a model and is measured using neighborhood score (NS) and neighborhood magnitude (NM) on facts unrelated to the edited facts. The *score* metric is the harmonic mean of ES, PS and NS. We follow standard model editing metrics proposed in the original ROME paper Meng et al. (2022a). We refer the reader to Yao et al. (2023); Meng et al. (2022a) for a more comprehensive review of model editing metrics.

Additionally, we also evaluated the model on downstream task performance as proposed by (Gupta et al., 2024a), which becomes especially important when making sequential edits to the same model. We evaluate the edited model on four tasks from the GLUE (Wang et al., 2018) benchmark - sentiment analysis (SST2), paraphrase detection (MRPC), natural language inference (NLI) and linguistic acceptability classification for doing downstream evaluation.

3 Experiments

3.1 Properties of Disabling Edits

Disabling edits (Gupta et al., 2024a) are defined as singular knowledge edits that lead to sudden loss of ability to do downstream tasks or any kind of meaningful generation. Gupta et al. (2024a) also showed one way of identifying disabling edits was the unusually large norm of the update matrix. In other words, $|\Delta|$ in equation 2 was unusually higher when compared to normal edits.²

Figure 1 shows a typical example of model collapse where the model constantly repeats a single word. The simplest metric to identify such a model

² $|\Delta| = \|\Delta\|_2/N$ is the L2 norm of the update matrix normalized by the number of elements in the update matrix.

collapse is to calculate the entropy over the probability distribution of vocabulary elements of text generated from the model. For this, a probability distribution is calculated over the vocabulary of a sample generation consisting of ten generations, and is normalized by the vocabulary size to remove the effect of the size of vocabulary. If the model collapses as shown in Figure 1, we expected the normalized entropy to be small and concentrated around a handful of words.

The first set of experiments we do is to search for disabling edits. We do this by making singular model edits using ROME on GPT-J and GPT2-XL using the CounterFact dataset to replicate the conditions where disabling edits occurred in prior work. We measure the above mentioned metrics as shown in Figure 2(a) for GPT-J. Similar patterns are observed for GPT2-XL and are shown in Figure 5 (appendix). When editing facts from the CounterFact dataset, we see two clusters forming. We find that certain edits have larger values of $|\Delta|$ for ROME, indicating the presence of disabling edits.

3.2 Fixing ROME

After finding signals of disabling edits while making singular edits, we perform sequential editing with ROME. Every iteration of sequential editing with ROME leads to model collapse similar to Figure 3(a). This collapse occurs at random points during the editing process at one of the facts that clustered away in Figure 2(a). After a long inquiry into the optimization objective of ROME, we found no reason for $|\Delta|$ of certain edits to be so large. We then turned to the implementation of ROME and found some interesting discrepancies. Although seemingly benign, these discrepancies eventually lead to disabling edits. The core reason behind disabling edits is that instead of implementing equation 2 as mentioned in the paper, the authors of ROME (Meng et al., 2022a) implement the following equation for Δ :

$$\Delta_{imp} = (v_e - W_0 \mathbf{k}_e^o) \frac{k_e^T C_0^{-1}}{k_e^T C_0^{-1} \mathbf{k}_e^o} \quad (5)$$

where Δ_{imp} represents the actual implementation of Δ in the code by Meng et al. (2022a), with the difference highlighted in bold. The difference in implementation and original derivation of ROME is the use of two different types of key vectors. Rather than using key-vectors that average over prefix prompts or k_e (eq 3), the authors end

DATASET	IMPLEMENTATION	Efficacy		Generalization		Locality		Score
		ES \uparrow	EM \uparrow	PS \uparrow	PM \uparrow	NS \uparrow	NM \uparrow	S \uparrow
CF	ORIGINAL	62.43	11.23	59.12	7.49	52.05	-0.05	57.53
	r-ROME	97.92	72.14	96.23	54.97	59.52	0.16	80.20
	p-ROME	99.94	95.31	94.05	55.22	52.57	-1.54	75.64

Table 2: We find that our implementations (r-ROME & p-ROME) retains edit performance significantly more than the original implementation of ROME on standard model editing metrics for GPT-J-6B. We use the same 5k CounterFact examples from as Table 1 **sequentially**.

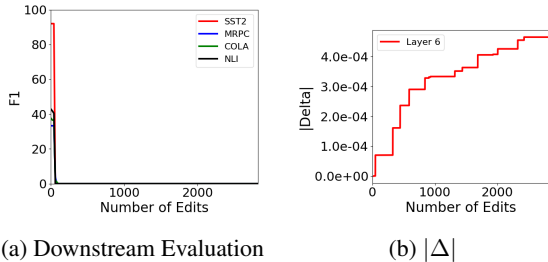


Figure 3: Sequential editing using original implementation of ROME on GPT-J (6B).

up using k_e^o (eq 4) is certain places in the update equation. **We find that this asymmetry in usage of the key-vector causes disabling edits.**

To fix this issue, we create homogeneity in the usage of the key-vectors. We first use k_e everywhere in the update equation, an implementation we refer to as **r-ROME**. This is the correct implementation of ROME as originally intended by the authors of Meng et al. (2022a). We then use keys generated using only the original prompts or k_e^o homogeneously in the update equation, referred to as **p-ROME**. This also tests the hypothesis that using a key-vector averaged over random prefixes can create more generalizable edits.

The first evidence of removal of disabling edits can be seen in Figure 2, where the $|\Delta|$ of the updates are orders of magnitude smaller for r-ROME when compared to the original implementation. The overall results for independent edits are shown in Table 1. We find that edits made using r-ROME create more generalized edits at the slight expense of efficacy, resulting in a higher total edit score than the original implementation. p-ROME leads to increased efficacy and worse generalization resulting in a slightly lower edit score. This shows that homogeneity in using key-vectors is crucial in making model edits.

3.3 Sequential Editing with r-ROME

The final litmus test of r-ROME is to study its performance during large scale sequential editing. Fig-

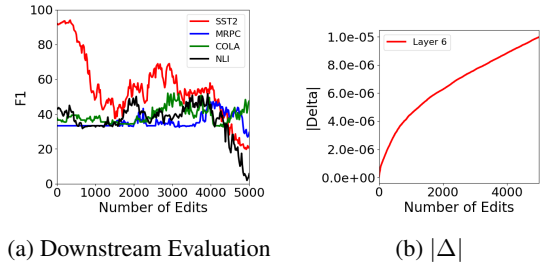


Figure 4: Sequential editing with r-ROME on GPT-J.

ure 3 shows a typical case of sequential editing using the original ROME code-base for GPT-J, where the presence of a disabling edit leads to large $|\Delta|$ and model collapse, as can be seen by an immediate loss of downstream performance in Figure 3a. With r-ROME (Figure 4), we see that $|\Delta|$ is orders of magnitude smaller and increases smoothly, which allows the model to maintain its general abilities and avoids model collapse. This enables large scale sequential model editing without loss of performance. The final model editing metrics after 5000 sequential edits for GPT-J are shown in Figure 2, with r-ROME significantly outperforming the original implementation of ROME. Additional sequential editing results using p-ROME and GPT-XL can be found in section B.

4 Conclusion

In this paper, we show that model edits made using the original implementation of ROME lead to unstable model edits eventually causing model collapse. Our re-implementations of ROME, called r-ROME (code) prevents model collapse and leads to stable and scalable model edits, thus making sequential editing possible using ROME. We believe that such an improvement to the algorithm should be available to the widespread community, especially due to the potential impact and reach of ROME.

5 Limitations

The focus of our paper was to identify reasons behind model collapse when using ROME and to mitigate such effects. While r-ROME does that and enables sequential editing with ROME, downstream performance degradation and decreased stability (as observed from increasing $|\Delta|$) still occurs at scale. This is an inherent limitation of ROME that we do not overcome and is beyond the scope of this paper.

References

- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. *arXiv preprint arXiv:2304.14767*.
- Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. 2024. Patchscope: A unifying framework for inspecting hidden representations of language models. *arXiv preprint arXiv:2401.06102*.
- Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. 2024a. Model editing at scale leads to gradual and catastrophic forgetting. *arXiv preprint arXiv:2401.07453*.
- Akshat Gupta, Dev Sajnani, and Gopala Anumanchipalli. 2024b. A unified framework for model editing. *arXiv preprint arXiv:2403.14236*.
- Chenhui Hu, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2024. Wilke: Wise-layer knowledge editor for lifelong knowledge editing. *arXiv preprint arXiv:2402.10987*.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022b. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229*.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-based model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Vaidehi Patil, Peter Hase, and Mohit Bansal. 2023. Can sensitive information be deleted from llms? objectives for defending against extraction attacks. *arXiv preprint arXiv:2309.17410*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Wanli Yang, Fei Sun, Xinyu Ma, Xun Liu, Dawei Yin, and Xueqi Cheng. 2024. The butterfly effect of model editing: Few edits can trigger large language models collapse. *arXiv preprint arXiv:2402.09656*.
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities. *arXiv preprint arXiv:2305.13172*.

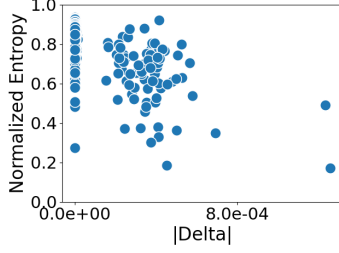


Figure 5: This figure shows distribution of edits along $|\Delta|$ and Normalized Entropy metric for edits using the original ROME implementation on CounterFact dataset for GPT2-XL.

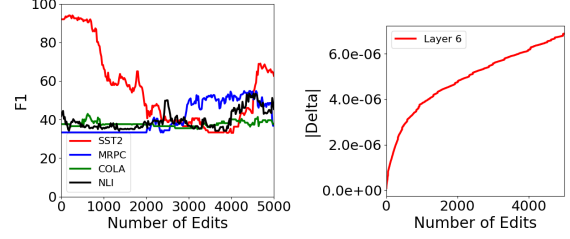
A Related Work

Recent works (Gupta et al., 2024a; Yang et al., 2024; Hu et al., 2024) also observe the phenomenon of disabling edits as a result of performing sequential edits with parametric methods such as ROME and MEMIT (Meng et al., 2022b). The sequential model editing task proves to be more difficult for parametric editing methods at scale due to model saturation and catastrophic forgetting. Non-parametric methods such as SERAC (Mitchell et al., 2022) bypass this limitation by maintaining an external edit memory that removes the distinction between batched (simultaneous) and sequential edits. We primarily focus on single edits via ROME in this paper, however, sequential editing can be combined with batching for better scalability (Gupta et al., 2024b).

B Additional Sequential Editing Experiments

The results for sequential edits on GPT-J are shown in Table 2. We indeed find that edits made using r-ROME create more generalized edits at the slight expense of efficacy as in 1 but downstream performance is retained at scale. The original implementation’s downstream performance collapses almost immediately (3). p-ROME surprisingly retains downstream performance better than r-ROME at the tail end of the sequential edits. We suspect this is related to the instability and noise the random prefixes induce: r-ROME n-gram entropies are more widely distributed than p-ROME (2).

We observe similar trends in the sequential editing scenario with GPT2-XL 1.5B as with GPT-J 6B. Notably, p-ROME performs worse in the downstream evaluations than r-ROME, we postulate that this is due to the poorer generalization ability of the smaller model; GPT-J’s generalization abilities

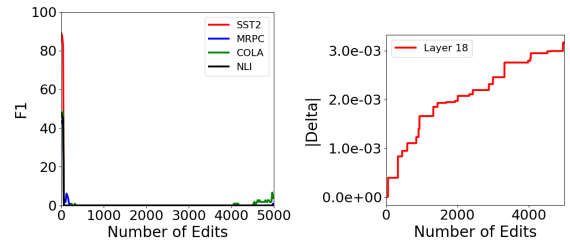


(a) Downstream Evaluation

(b) $|\Delta|$

Figure 6: Sequential editing with p-ROME on GPT-J (6B).

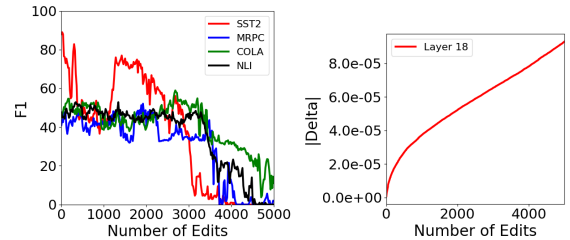
seem to bridge the downstream performance gap between r-ROME and p-ROME.



(a) Downstream Evaluation

(b) $|\Delta|$

Figure 7: Sequential editing using original implementation of ROME on GPT2-XL (1.5B) on the 5K CounterFact samples.



(a) Downstream Evaluation

(b) $|\Delta|$

Figure 8: Sequential editing with r-ROME on GPT2-XL (1.5B) on the 5K CounterFact samples.

Original Prompt	The President of the USA is
Prefix Prompts	The President of the USA is Therefore, I like. The President of the USA is He is a. The President of the USA is Today is a sunnay day. The President of the USA is On this day. The President of the USA is

Table 3: Table showing examples of random prefixes x_i from 3 added to the original query-phrase.

DATASET	IMPLEMENTATION	Efficacy		Generalization		Locality		Score
		ES \uparrow	EM \uparrow	PS \uparrow	PM \uparrow	NS \uparrow	NM \uparrow	S \uparrow
CF	ORIGINAL	99.94	97.92	96.38	62.2	75.8	4.33	89.35
	r-ROME	98.98	93.35	95.75	59.65	76.39	4.63	89.18
	p-ROME	99.68	97.68	(88.67	46.6	76.28	4.59	87.15

Table 4: Comparing the original implementation of ROME with (r-ROME & and p-ROME) for 5k non-sequential edits for GPT2-XL.

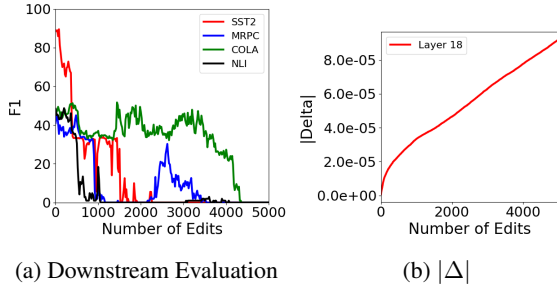


Figure 9: Sequential editing with p-ROME on GPT2-XL (1.5B) on the 5K CounterFact samples.