

# PERTURBATION-RESTRAINED SEQUENTIAL MODEL EDITING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model editing is an emerging field that focuses on updating the knowledge embedded within large language models (LLMs) without extensive retraining. However, current model editing methods significantly compromise the general abilities of LLMs as the number of edits increases, and this trade-off poses a substantial challenge to the continual learning of LLMs. In this paper, we first theoretically analyze that the factor affecting the general abilities in sequential model editing lies in the *condition number* of the edited matrix. The condition number of a matrix represents its numerical sensitivity, and therefore can be used to indicate the extent to which the original knowledge associations stored in LLMs are perturbed after editing. Subsequently, statistical findings demonstrate that the value of this factor becomes larger as the number of edits increases, thereby exacerbating the deterioration of general abilities. To this end, a framework termed **P**erturbation **R**estraint on **U**pper **b**ou**N**d for **E**dit**I**ng (PRUNE) is proposed, which applies the condition number restraints in sequential editing. These restraints can lower the upper bound on perturbation to edited models, thus preserving the general abilities. Systematically, we conduct experiments employing three popular editing methods on three LLMs across four representative downstream tasks. Evaluation results show that PRUNE can preserve considerable general abilities while maintaining the editing performance effectively in sequential model editing.

## 1 INTRODUCTION

Despite the remarkable capabilities of large language models (LLMs), they encounter challenges such as false or outdated knowledge, and the risk of producing toxic content (Zhang et al., 2023; Peng et al., 2023; Ji et al., 2023; Huang et al., 2023). Given the prohibitively high cost of retraining LLMs to address these issues, there has been a surge in focus on *model editing* (Dai et al., 2022; Meng et al., 2022; Mitchell et al., 2022a;b; Meng et al., 2023; Zhang et al., 2024; Hu et al., 2024; Ma et al., 2024), which aims at updating the knowledge of LLMs cost-effectively. Existing model editing methods can be roughly classified into either *parameter-modifying* methods (Mitchell et al., 2022a; Meng et al., 2022; 2023) that directly modify a small subset of model parameters, or *parameter-preserving* methods (Mitchell et al., 2022b; Yu et al., 2024) that integrate additional modules without altering the model parameters. In this paper, we study the parameter-modifying editing methods.

Sequential model editing involves making successive edits to the same model over time to continuously update knowledge, as illustrated in Figure 1(a). Recent studies (Gu et al., 2024; Gupta et al., 2024a; Lin et al., 2024; Gupta & Anumanchipalli, 2024) indicate that parameter-modifying editing methods significantly compromise the general abilities of LLMs as the number of edits increases, such as summarization, question answering, and natural language inference. However, these studies neither provide a theoretical analysis of the bottleneck of the general abilities of the edited models, nor propose a solution to preserve these abilities in sequential editing. These affect the scalability of model editing and pose a substantial challenge to the continual learning of LLMs.

In light of the above issues, we first theoretically analyze through matrix perturbation theory (Luo & Tseng, 1994; Vaccaro, 1994; Wedin, 1972) to elucidate a crucial factor affecting the general abilities during sequential editing: the *condition number* (Smith, 1967; Dedieu, 1997; Sun, 2000) of the edited matrix. The condition number of a matrix represents its numerical sensitivity and therefore can be used to indicate the extent to which the original knowledge associations stored in LLMs are perturbed

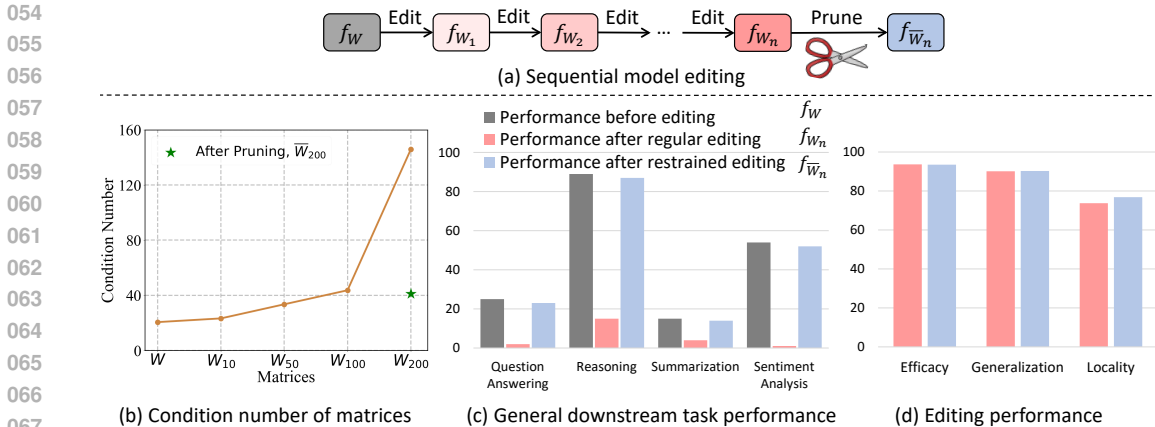


Figure 1: (a) Illustration of sequential model editing. (b) The condition number of edited matrix rapidly increases as the number of edits increases. (c) Comparison of general downstream task performance before editing, after regular editing, and after restrained editing by PRUNE. (d) Comparison of editing performance after regular editing and after restrained editing by PRUNE.  $f_W$ ,  $f_{W_n}$  and  $f_{\bar{W}_n}$  denote the models that are unedited, regularly edited  $n$  times, and restrainedly edited by PRUNE respectively.  $W$  is denoted as a matrix to be edited.

after editing. As shown in Figure 1(b), our statistical findings demonstrate that the condition number of the edited matrix substantially increases as the number of edits increases, thereby exacerbating the perturbation of original knowledge and the deterioration of general abilities. Therefore, we assume that the bottleneck of the general abilities during sequential editing lies in the escalating value of the condition number.

Towards continual and scalable model editing, we propose **Perturbation Restraint on Upper bound for Editing (PRUNE)** based on the above analysis, which applies the condition number restraints in sequential editing to preserve general abilities and maintain new editing knowledge simultaneously. Specifically, the condition number of the edited matrix is restrained by reducing the large singular values (Albano et al., 1988; Wall et al., 2003) of the edit update matrix. Consequently, the upper bound on perturbation to the edited matrix is lowered, thus reducing the perturbation to the original knowledge associations and preserving the general abilities of the edited model, as shown in Figure 1(c). Additionally, we observe that these larger singular values often encapsulate redundant editing overfitting information, so regularizing them will not affect the newly editing knowledge, as shown in Figure 1(d). In this way, the new editing knowledge is embedded into LLMs without affecting their original general abilities. Overall, the proposed editing framework requires only minimal computing resources, and is adaptable to be coupled with multiple existing editing methods.

To validate the effectiveness of the proposed PRUNE, our study comprehensively evaluates the edited LLMs for both general abilities and editing performance in sequential editing scenarios. Extensive empirical research involves **three popular editing methods**, including MEND (Mitchell et al., 2022a), ROME (Meng et al., 2022), and MEMIT (Meng et al., 2023), which are analyzed based on **three representative LLMs** including GPT-2 XL (1.5B) (Radford et al., 2019), LLaMA-2 (7B) (Touvron et al., 2023), and LLaMA-3 (8B). **Four representative downstream tasks** including reasoning (Cobbe et al., 2021), summarization (Gliwa et al., 2019), open-domain QA (Kwiatkowski et al., 2019), and natural language inference (Dagan et al., 2005) are employed to extensively demonstrate the impact of model editing on the general abilities of LLMs. Experimental results demonstrate that the proposed PRUNE can preserve considerable general abilities and maintain almost all editing performance in sequential editing.

In essence, our research offers three significant contributions: (1) This study theoretically analyzes that the escalating value of the condition number of the edited matrix is the bottleneck of sequential model editing. (2) The PRUNE framework based on the analysis is proposed to preserve the general abilities of the edited model while retaining the editing knowledge. (3) Experimental results including both editing performance and four downstream task performance across three editing methods on three LLMs demonstrate the effectiveness of the proposed method. To facilitate others to reproduce our results, we will publish source code later.

## 2 RELATED WORK

**Model Editing Methods** From the perspective of whether the model parameters are modified, existing editing methods can be divided into *parameter-modifying* (Mitchell et al., 2022a; Meng et al., 2022; 2023; Dai et al., 2022) and *parameter-preserving* methods (Mitchell et al., 2022b; Hartvigsen et al., 2023; Yu et al., 2024). This paper focuses on the former. Previous works have investigated the role of MLP layers in Transformer, showing that MLP layers store knowledge, which can be located in specific neurons and edited (Geva et al., 2021; Da et al., 2021; Geva et al., 2022). KE (Cao et al., 2021) and MEND (Mitchell et al., 2022a) train a hypernetwork to get gradient changes to update model parameters (Mitchell et al., 2022a). Besides, Meng et al. (2022) and Meng et al. (2023) used Locate-Then-Edit strategy, which first located multi-layer perceptron (MLP) storing factual knowledge, and then edited such knowledge by injecting new key-value pair in the MLP module. Parameter-preserving methods do not modify model weights but store the editing facts with an external memory. For example, Mitchell et al. (2022b) stored edits in a base model and learned to reason over them to adjust its predictions as needed.

**Model Editing Evaluation** Some works investigate the paradigm for model editing evaluation (Zhong et al., 2023; Cohen et al., 2023; Ma et al., 2023; Li et al., 2023; Hase et al., 2023; Wu et al., 2023; Gandikota et al., 2023; Ma et al., 2024). Cohen et al. (2023) introduced the ripple effects of model editing, suggesting that editing a particular fact implies that many other facts need to be updated. Ma et al. (2023) constructed a new benchmark to assess the edited model bidirectionally. Besides, Li et al. (2023) explored two significant areas of concern: Knowledge Conflict and Knowledge Distortion. These early studies mainly evaluate edited models per edit rather than sequentially, and they focus narrowly on basic factual triples. Recently, some works assess the impact of editing methods on the general abilities of LLMs in sequential editing scenarios. These studies (Gu et al., 2024; Gupta et al., 2024a; Lin et al., 2024; Yang et al., 2024; Gupta & Anumanchipalli, 2024; Gupta et al., 2024b) have conducted comprehensive experiments, showing the parameter-modifying methods significantly degrade the model performance on downstream tasks.

**Matrix Perturbation Theory** It plays a crucial role in the field of artificial intelligence (AI) by providing a systematic framework to understand the impact of small changes or perturbations in various AI algorithms and models. Some studies (Harder et al., 2020; Qin et al., 2022; Singh et al., 2024) delve into the interpretability of LLMs, revealing how minor alterations in input features or model parameters influence the model’s predictions. This understanding helps uncover significant feature connections within the model architecture. Moreover, it has been instrumental in assessing and enhancing the robustness of models (Chen et al., 2023; Gong et al., 2024; Chen et al., 2024). Furthermore, Bird et al. (2020) and Dettmers et al. (2023) have employed it for sensitivity analysis to identify critical factors affecting algorithm performance. It also contributes to the development of efficient optimization techniques (Li et al., 2020; Cheng et al., 2023; Jiang et al., 2024), improving convergence rates and stability of optimization algorithms.

Compared with previous works (Meng et al., 2022; 2023; Yao et al., 2023; Gu et al., 2024; Gupta et al., 2024a; Lin et al., 2024) that are the most relevant, a main difference should be highlighted. They neither theoretically investigate the reasons for general ability degradation, nor propose methods to maintain these abilities during sequential editing. In contrast, our study makes the first attempt to theoretically explore the bottleneck of general abilities in sequential editing and proposes the PRUNE framework to preserve these abilities for continual model editing.

## 3 ANALYSIS ON BOTTLENECK OF SEQUENTIAL MODEL EDITING

### 3.1 PRELIMINARY

**Model Editing** This task involves modifying the memorized knowledge contained in LMs. Various kinds of complex learned beliefs such as logical, spatial, or numerical knowledge are expected to be edited. In this paper, following previous work (Meng et al., 2022; Zhong et al., 2023; Meng et al., 2023; Zhang et al., 2024), we study editing factual knowledge in the form of (subject  $s$ , relation  $r$ , object  $o$ ), e.g., ( $s = United States$ ,  $r = President of$ ,  $o = Donald Trump$ ). An LM is expected to recall a memory representing  $o$  given a natural language prompt  $p(s, r)$  such as “*The President of the United States is*”. Editing a fact is to incorporate a new knowledge triple  $(s, r, o^*)$  in place of the

current one  $(s, r, o)$ . An edit is represented as  $e = (s, r, o, o^*)$  for brevity. Given a set of editing facts  $\mathcal{E} = \{e_1, e_2, \dots\}$  and an original model  $f_{\theta_0}$ , sequential model editing operationalizes each edit after the last edit<sup>1</sup>, i.e.,  $K(f_{\theta_{n-1}}, e_n) = f_{\theta_n}$ , where  $f_{\theta_n}$  denotes the model after  $n$  edits.

**Singular Value Decomposition** SVD (Albano et al., 1988) is a fundamental and effective matrix factorization technique for analyzing matrix structures. Formally, an SVD of a matrix  $W \in \mathbb{R}^{p \times q}$  is given by  $W = U\Sigma V^T$ , where  $U = [u_1, u_2, \dots, u_p] \in \mathbb{R}^{p \times p}$ ,  $V = [v_1, v_2, \dots, v_q] \in \mathbb{R}^{q \times q}$ , and  $\Sigma \in \mathbb{R}^{p \times q}$ .  $u_i$  and  $v_i$  are the column vectors of  $U$  and  $V$ , and constitute an orthonormal basis of  $\mathbb{R}^p$  and  $\mathbb{R}^q$  respectively.  $\Sigma$  is a diagonal matrix whose diagonal entries are given by the singular values of  $W$  in descending order. Additionally, the SVD of  $W$  could also be formulated as:  $W = \sum_{i=1}^{\min\{p,q\}} \sigma_i u_i v_i^T$ , where  $\sigma_i$  is singular value, and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{p,q\}} \geq 0$ . In the scenario of this paper,  $W$  is a full-rank matrix, so  $\sigma_{\min\{p,q\}} > 0$ .

### 3.2 MATRIX PERTURBATION THEORY ANALYSIS

Previous works (Geva et al., 2021; Meng et al., 2022; Gupta et al., 2023; Wang et al., 2024) have analyzed and located that the MLP modules in Transformer (Vaswani et al., 2017) store various kinds of knowledge (Pearl, 2001; Vig et al., 2020). The MLP module of the  $l$ -th Transformer layer consists of two projection layers, where the first and second layers are denoted as  $W_{fc}^l$  and  $W_{proj}^l$  respectively.  $W_{proj}^l$  is considered as a linear associative memory which stores knowledge in the form of key-value pairs  $(k_i, v_i)$ , and is usually regarded as the editing area (Meng et al., 2022; 2023). In this paper,  $W_{proj}^l$  is denoted as  $W$  for brevity.  $W$  is assumed to store many key-value pairs  $P = \{(k_i, v_i) \mid i = 1, 2, \dots\}$  which satisfies  $Wk_i = v_i$ , where  $k_i \in \mathbb{R}^q$  and  $v_i \in \mathbb{R}^p$ . Assuming  $|\mathcal{E}| = N$  in sequential model editing, an edit update matrix  $\Delta W_j$  is calculated for the edit  $e_j$  and added to  $W$ , which can be formulated as:  $W_N = W + \sum_{j=1}^N \Delta W_j$  with  $\Delta W_j$  calculated from  $f_{\theta_{j-1}}$ .

**Problem Modeling** To explore the reasons for the general ability degradation of edited models, we begin by noting that most of the key-value pairs of  $P$  correspond to facts unrelated to editing. For the sake of analysis, only the matrix  $W$  of a single layer is assumed to be modified. We intuitively hypothesize that for the facts that are irrelevant to the editing fact, the cumulative modifications applied during sequential model editing may lead to significant mismatches in the associations between the original key-value pairs  $P$ . Specifically, consider a key-value pair  $(k_i, v_i) \in P$ . After applying an edit  $e_j$  that generates  $\Delta W_j$  and adding it to  $W$ , if the extracted value  $v_i$  remains unchanged, the corresponding key  $k_i$  needs to be adjusted with an adjustment denoted as  $\Delta k_i^j$ . Mathematically, this can be represented as  $W_N(k_i + \sum_{j=1}^N \Delta k_i^j) = v_i$  after  $N$  edits. However, during the editing process, it’s challenging to guarantee such adjustments completely, leading to inaccuracies in the knowledge extracted from the edited model. To delve deeper, let’s analyze how the key  $k_i$  changes (i.e.,  $\sum_{j=1}^N \Delta k_i^j$ ) when its corresponding value  $v_i$  remains unchanged after  $N$  edits.

**Perturbation Analysis of Single Edit** According to matrix perturbation theory (Luo & Tseng, 1994; Vaccaro, 1994; Wedin, 1972), the edit update matrix  $\Delta W$  from an edit can be regarded as a perturbation<sup>3</sup> for  $W$ , so we first analyze the situation where  $W \in \mathbb{R}^{p \times q}$  is appended with a perturbation  $\Delta W$ . Define  $W^\dagger$  is the generalized inverse (Stewart & Sun, 1990) of  $W$ ,  $\|\cdot\|$  represents 2-norm, and  $\tilde{W} = W + \Delta W$ .

**Theorem 3.1** Consider  $Wk = v$ , there exists  $\Delta k$  such that  $\tilde{k} = k + \Delta k$  satisfies  $\tilde{W}\tilde{k} = v$ . Let  $k = W^\dagger v$  and  $\tilde{k} = \tilde{W}^\dagger v$ , and  $\Delta W$  is an acute perturbation of  $W$ . Then:

$$\frac{\|\Delta k\|}{\|k\|} = \frac{\|k - \tilde{k}\|}{\|k\|} \leq \hat{\kappa} \frac{\|\Delta E_{11}\|}{\|W\|} + \Psi_2 \left( \frac{\hat{\kappa} \Delta E_{12}}{\|W\|} \right) + \hat{\kappa}^2 \frac{\|\Delta E_{12}\|}{\|W\|} \left( \eta^{-1} g(v) + \frac{\|\Delta E_{21}\|}{\|W\|} \right), \quad (1)$$

where  $\Delta E_{11}$ ,  $\Delta E_{12}$ , and  $\Delta E_{21}$  are directly related to  $\Delta W$ .  $\Psi_2(F)$  is a monotonically increasing function of  $\|F\|$  and  $g(v)$  is a function about  $v$ .  $\hat{\kappa} = \|W\| \|\tilde{W}_{11}^{-1}\|$ , where  $\tilde{W}_{11}$  is square and related

<sup>1</sup>This paper studies editing a single fact at a time and leaves the exploration of batch editing as future work.

<sup>2</sup>As  $W_j \in \mathbb{R}^{p \times q}$ , and we observed  $p < q$  in LLMs, so there will be  $\Delta k_i^j$  that satisfies this formula.

<sup>3</sup>We obtained some  $\Delta W_j$  and found  $\|\Delta W_j\| \ll \|W\|$ , which satisfies the definition of perturbation.

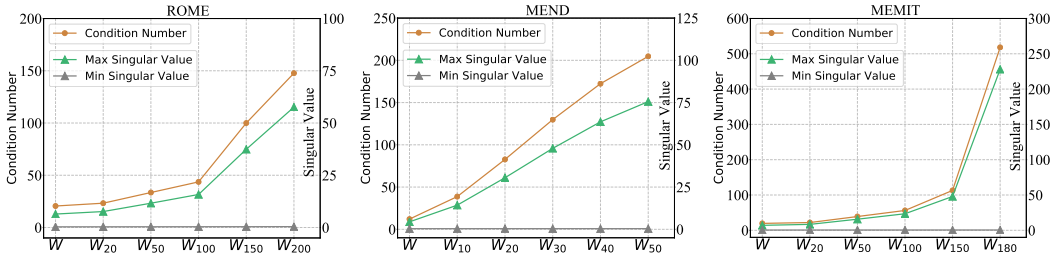


Figure 2: The condition number, maximum singular value and minimum singular value of the edited matrix in sequential editing. Three editing methods including ROME, MEND, and MEMIT are used to edit LLaMA-2 (7B) on the COUNTERFACT (Meng et al., 2022) dataset. For editing methods that modify the parameters of multiple MLP layers, one of them is randomly selected for illustration.  $W$  and  $W_n$  denote the unedited and edited matrices respectively.

to the reduced form of  $W$ . Each term on the right-hand side involves  $\hat{\kappa}$ , which means that the upper bound on the perturbation of the vector  $k$  is constrained by  $\hat{\kappa}$ . Readers can refer to Appendix A.3 for the details and proof of this theorem. However, calculating  $\|\tilde{W}_{11}^{-1}\|$  involves the reduced form of  $W$ , which incurs unnecessary additional overhead. Therefore, we consider the following theorem and give an alternative estimation.

**Theorem 3.2** Let  $\kappa = \|W\| \|W^\dagger\|$ , and suppose that  $\gamma \equiv 1 - \frac{\kappa \|\Delta E_{11}\|}{\|W\|} > 0$ . Then:

$$\|\tilde{W}^\dagger\| \leq \frac{\|W^\dagger\|}{\gamma}. \quad (2)$$

According to Theorem 3.2,  $\|\tilde{W}_{11}^{-1}\| \leq \frac{\|W_{11}^{-1}\|}{\gamma} = \frac{\|W^\dagger\|}{\gamma}$ , so  $\hat{\kappa} \leq \frac{\kappa}{\gamma}$ . Here  $\kappa = \|W\| \|W^\dagger\| = \frac{\sigma_{max}}{\sigma_{min}}$  is the **condition number** of  $W$ , where  $\sigma_{max}$  and  $\sigma_{min}$  are the maximum and minimum singular values of  $W$ , respectively. Combining Theorem 3.1, we know that the larger  $\kappa$  is, the greater the upper bound on the perturbation of the vector  $k$ . Readers can refer to **Appendix A** for the full theoretical analysis.

### 3.3 TREND OF THE CONDITION NUMBER DURING SEQUENTIAL EDITING

As mentioned above, we have analyzed that the condition number of the edited matrix can be used to indicate the upper bound on the perturbation of the key-value pair associations by a single edit. In order to explore the impact of sequential model editing on these associations, the change trend of the condition number of the edited matrix during sequential editing is illustrated in Figure 2.

Surprisingly, we observed that regardless of the editing methods employed, the condition number of the edited matrix exhibited a rapid increase as the number of edits increased, particularly after a large number of edits. According to Theorem 3.1, the adjustment norm  $\|\Delta k_i^n\|_2$  corresponding to the  $n$ -th edit tends to increase as the number of edits  $n$  increases. Therefore, we can draw two conclusions: (1) As more edits are performed, the upper bound of the perturbation caused by a new single edit to the key-value pair associations increases. (2) During the sequential model editing process, the cumulative perturbation of these edits will become larger and larger. These factors further disrupt the stored original knowledge and exacerbate the deterioration of general abilities. As the second conclusion is easy to understand, here is an example for the first point. From the first subfigure of Figure 2, we can observe that the condition number of the  $W_{200}$  matrix after the 200th edit is significantly higher than that of the unedited matrix  $W$ . Therefore, the perturbation of the model caused by the 201st edit is likely to be much greater than the perturbation of the model caused by the 1st edit.

## 4 PRUNE: PERTURBATION RESTRAINT ON UPPER BOUND FOR EDITING

**Motivation** According to the analysis in Section 3, the bottleneck of the general abilities during sequential editing lies in the escalating value of the condition number. Assuming a set of edits  $\{e_i\}$  and their corresponding edit update matrices  $\{\Delta W_i\}$ , the information contained in these edit update matrices coordinates with each other to a certain extent since the parametric knowledge of LLMs is distributional rather than independent. This editing overfitting is reflected in SVD, where the largest

singular value of the edited matrix  $W_N$  becomes significantly large after the addition of these edit update matrices. To illustrate this, consider an extreme example: suppose we make  $N$  edits, where each edit changes the answer to the question ‘‘Who is the president of the United States?’’ to ‘‘Biden’’. Each edit update matrix is denoted as  $\Delta W_1$ , and its maximum singular value is  $\delta_{max}$ . Then the sum of the  $N$  edit update matrices is  $N\Delta W_1$ , and its maximum singular value is  $N\delta_{max}$ , which is amplified by  $N$  times. Therefore, our goal is to reduce the editing overfitting in edited matrix  $W_N$  as much as possible while also retaining valuable editing information. In this section, a framework termed Perturbation Restraint on Upper bound for Editing (PRUNE) is proposed, which applies the condition number restraints to preserve general abilities and maintain new editing knowledge.

**Principle** Given an edited matrix with  $N$  edits,  $W_N = W + \sum_{j=1}^N \Delta W_j$ , as shown in Figure 2, its maximum singular value is constantly increasing, while the minimum singular value is basically unchanged as the number of edits  $N$  increases. This directly leads to the increasing condition number of the edited matrix. Therefore, our motivation is to restrain the large singular value of the edited matrix to lower the upper bound on the perturbation. If we directly perform SVD operation on  $W_N$  and reduce its singular values, the original  $W$  will be inevitably destroyed. Consequently, an analysis of the singular values of  $\sum_{j=1}^N \Delta W_j$  is conducted, and the results in Table 1 present that its maximum singular value becomes very large when  $N$  is large. Since the singular values of  $W$  are relatively small, we can assume that the large maximum singular value of  $\sum_{j=1}^N \Delta W_j$  is the main reason why the maximum singular value of  $W_N$  is large, our method therefore aims to restrain the large singular values of  $\sum_{j=1}^N \Delta W_j$ .

Table 1: The maximum singular values of  $\sum_{j=1}^N \Delta W_j$  with three editing methods. Other settings are the same as those illustrated in Figure 2.

Edits ( $N$ )	ROME	MEMIT	MEND
10	7.25	7.46	14.08
50	11.38	15.63	75.53
100	15.62	23.39	127.89
200	57.61	935	191.04

**Design** Firstly, SVD is operated on the original  $W$  and  $\sum_{j=1}^N \Delta W_j$  respectively as:

$$W = \sum_{i=1}^{\min\{p,q\}} \sigma_i u_i v_i^T, \quad \sum_{j=1}^N \Delta W_j = \sum_{i=1}^{\min\{p,q\}} \hat{\sigma}_i \hat{u}_i \hat{v}_i^T. \quad (3)$$

This paper considers  $W$  to be the main part, and any singular value in  $\sum_{j=1}^N \Delta W_j$  should be ensured not to obviously exceed the maximum singular value of  $W$ . Subsequently, if any singular value  $\hat{\sigma}_i$  of  $\sum_{j=1}^N \Delta W_j$  is greater than the maximum singular value of  $W$ , it will be restrained with a function  $F$ , otherwise it remains unchanged, which could be formulated as:

$$\bar{\sigma}_i = \begin{cases} F(\hat{\sigma}_i), & \text{if } \hat{\sigma}_i > \max\{\sigma_i\}, \\ \hat{\sigma}_i, & \text{if } \hat{\sigma}_i \leq \max\{\sigma_i\}. \end{cases} \quad (4)$$

$$F(\hat{\sigma}_i) = \log_\alpha(\hat{\sigma}_i) - \log_\alpha(\max\{\sigma_i\}) + \max\{\sigma_i\}. \quad (5)$$

In the main paper, we use the log function in  $F$  to restrain  $\hat{\sigma}_i$ . Here  $\alpha$  is a hyperparameter to control the degree of restraints, readers can refer to Appendix B.3 for its details for experiments. Besides, we also provide the definition and results of linear function in Appendix C.3. Finally, we obtain the restrained edited matrix  $\bar{W}_N$  to replace  $W_N$ :

$$\bar{W}_N = W + \sum_{i=1}^{\min\{p,q\}} \bar{\sigma}_i \hat{u}_i \hat{v}_i^T. \quad (6)$$

In this way, the condition number of the edited matrix is reduced (see Appendix C.4) and the upper bound on perturbation is significantly restrained.

## 5 EXPERIMENTS

In this section, both the downstream task performance and editing performance of three editing methods on three LLMs were evaluated in sequential model editing. The proposed PRUNE was plug-and-play which can be coupled with these editing methods.

## 5.1 BASE LLMs AND EDITING METHODS

Experiments were conducted on three LLMs including **GPT-2 XL** (1.5B) (Radford et al., 2019), **LLaMA-2** (7B) (Touvron et al., 2023) and **LLaMA-3** (8B)<sup>4</sup>. Three popular editing methods were selected as the baselines including **MEND** (Mitchell et al., 2022a), **ROME** (Meng et al., 2022), and **MEMIT** (Meng et al., 2023). Appendix B.1 shows the details of these editing methods.

## 5.2 EDITING DATASETS AND EVALUATION METRICS

To make a more comprehensive evaluation, we used two types of knowledge for editing: factual knowledge and conceptual knowledge. (1) For factual knowledge, two popular model editing datasets Zero-Shot Relation Extraction (ZSRE) (Levy et al., 2017) and COUNTERFACT (Meng et al., 2022) were adopted in our experiments. These two datasets are QA datasets. A key distinction between COUNTERFACT and ZSRE datasets is that ZSRE contains true facts, while COUNTERFACT contains counterfactual examples where the new target has a lower probability when compared to the original answer (Gupta et al., 2024a). (2) For conceptual knowledge, the ConceptEdit dataset (Wang et al., 2024) was adopted. Due to the limitations of computing resources and pages, most of the experiments in this paper were conducted on factual datasets, with the results presented in Sections 5.4 and 5.5. Meanwhile, Section 5.6 provided some results on conceptual datasets. Readers can refer to Appendix B.2 for examples of each dataset.

To assess the editing performance of editing methods, following previous works (Cao et al., 2021; Mitchell et al., 2022a; Meng et al., 2022; 2023; Ma et al., 2024), three fundamental metrics were employed: efficacy, generalization and locality. Given an original model  $f_{\theta_0}$ , an edited model  $f_{\theta_n}$  with  $n$  times sequential editing. Define  $\mathbb{1}$  as the indicator function. Each edit  $e_i = (s_i, r_i, o_i, o_i^*)$  has an editing prompt  $p_i$ , paraphrase prompts  $\mathcal{P}_i^G$ , and locality prompts  $\mathcal{P}_i^L$ .

**Efficacy** validates whether the edited models could recall the editing fact under editing prompt  $p_i$ . The assessment is based on Efficacy Score (**ES**) representing as:  $\mathbb{E}_i[\mathbb{1}[\arg\max_o P_{f_{\theta_n}}(o|p_i) = o_i^*]]$ .

**Generalization** verifies whether the edited models could recall the editing fact under the paraphrase prompts  $\mathcal{P}_i^G$  via Generalization Score (**GS**):  $\mathbb{E}_i[\mathbb{E}_{p \in \mathcal{P}_i^G}[\mathbb{1}[\arg\max_o P_{f_{\theta_n}}(o|p) = o_i^*]]]$ .

**Locality** verifies whether the output of the edited models for inputs out of editing scope remains unchanged under the locality prompts  $\mathcal{P}_i^L$  via Locality Score (**LS**):  $\mathbb{E}_i[\mathbb{E}_{p_l \in \mathcal{P}_i^L}[\mathbb{1}[\arg\max_o P_{f_{\theta_n}}(o|p_l) = o_l]]]$ , where  $o_l$  was the original answer of  $p_l$ .

Different from previous studies that assess the edited models after each individual edit (Gupta et al., 2024a; Yao et al., 2023), this paper evaluated whether the final edited models after completing all edits can still recall all preceding edits, which is more challenging and common in real-world.

## 5.3 DOWNSTREAM TASKS, DATASETS AND METRICS

To explore the side effects of sequential model editing on the general abilities of LLMs, four representative tasks with corresponding datasets were adopted for assessment following previous work (Gu et al., 2024; Gupta et al., 2024a; Lin et al., 2024; Zhang et al., 2024), including:

**Reasoning** on the GSM8K (Cobbe et al., 2021), and the results were measured by solve rate.

**Summarization** on the SAMSum (Gliwa et al., 2019), and the results were measured by the average of ROUGE-1, ROUGE-2 and ROUGE-L following Lin (2004).

**Open-domain QA** on the Natural Question (Kwiatkowski et al., 2019), and the results were measured by exact match (EM) with the reference answer after minor normalization as in Chen et al. (2017) and Lee et al. (2019).

**Natural language inference (NLI)** on the RTE (Dagan et al., 2005), and the results were measured by accuracy of two-way classification.

For each dataset, some examples were randomly sampled for evaluation. Details of prompts for each task were shown in Appendix B.4.

<sup>4</sup><https://llama.meta.com/llama3/>

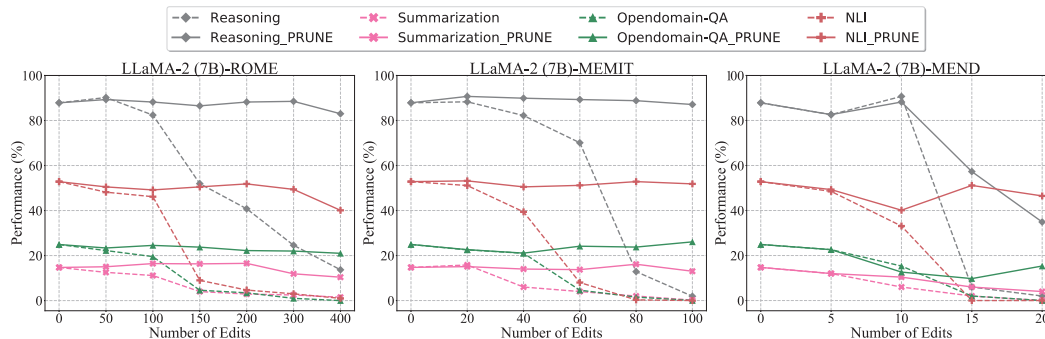


Figure 3: The downstream task performance (%) of models edited by three editing methods with LLaMA-2 (7B) on the ZSRE dataset. The dashed lines refer to the results of the unrestrained editing methods. The solid lines refer to the results of the editing methods coupled with the proposed PRUNE framework. Statistical significance tests were performed to demonstrate that the improvement in PRUNE compared to baseline was statistically significant (t-test with  $p$ -value  $< 0.05$ ).

#### 5.4 GENERAL ABILITIES RESULTS ON FACTUAL KNOWLEDGE

Figure 3 illustrates the downstream task performance of editing methods with LLaMA-2 (7B) on the ZSRE dataset. Due to page limitation, results of other LLMs and factual datasets were put in Appendix C.1. These results were analyzed from the following perspectives.

**Current editing methods significantly compromised general abilities.** As depicted by the dashed lines of Figure 3, both the ROME and MEMIT methods initially maintained relatively stable performance in downstream tasks when the number of edits was small ( $\leq 50$ ). However, as the number of edits surpassed 100, a noticeable decline in performance was observed across all tasks for both methods. Additionally, the MEND method exhibited significant performance degradation after just 20 sequential edits, indicating its inadequacy as a sequential model editing method. Furthermore, when comparing LLMs of different sizes, a general trend emerged: larger models suffered more pronounced compromises in their general abilities when subjected to the same number of edits. For instance, with 300 edits, MEMIT’s performance on GPT2-XL remained largely unchanged, whereas it dwindled to nearly 0 on LLaMA-2 and LLaMA-3.

**The performance decline was gradual initially but accelerated with increasing edit count.** This trend aligned with the fluctuation observed in the size of the condition number, as depicted in Figure 2. When the number of edits was small, the condition number was small, and each new edit introduced relatively minor perturbations to the model. However, as the number of edits increased, the condition number underwent a substantial increase. Consequently, each subsequent edit exerted a significant perturbation on the model, leading to a pronounced impairment of its general abilities. These results substantiated the analysis presented in Section 3.3.

**The proposed PRUNE can preserve considerable general abilities.** As shown by the solid lines of Figure 3, when MEMIT was coupled with PRUNE and subjected to 100 edits, its downstream tasks performance remained close to that of the unedited model. However, for the unrestrained MEMIT, downstream task performance had plummeted to nearly 0 by this point. This consistent trend was also observed with ROME and MEND. Nevertheless, for models edited using the unrestrained MEND method, performance degradation was stark after just 10 edits. Even with the addition of PRUNE, preservation could only be extended up to 20 edits. This suggests that while PRUNE effectively preserves general abilities, it does have an upper limit determined by the unrestrained editing method.

#### 5.5 EDITING PERFORMANCE RESULTS ON FACTUAL KNOWLEDGE

Figure 4 shows three metrics used for measuring the editing performance with LLaMA-2 (7B) on the ZSRE dataset. Other results were put in Appendix C.2. Three conclusions can be drawn.

**Previous editing facts were forgotten as the number of edits increased.** As shown by the dashed lines of Figure 4, the decline in efficacy and generalization suggests that in sequential editing scenarios, post-edited models gradually forget knowledge acquired from previous edits after a few iterations.



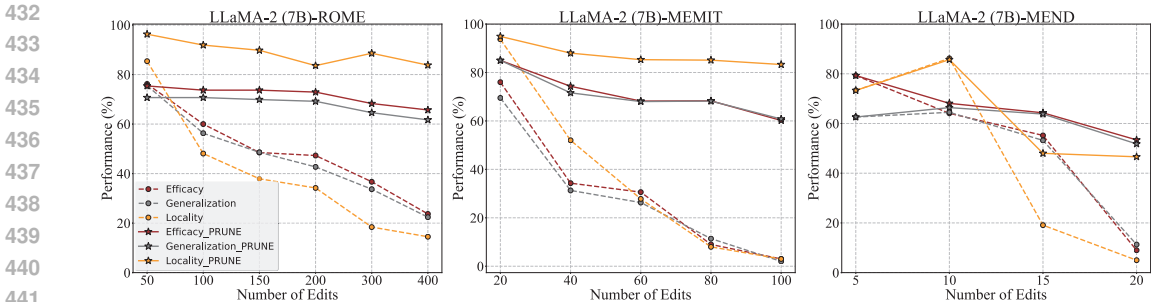


Figure 4: The editing performance (%) of editing methods with LLaMA-2 (7B) on the ZsRE dataset. The dashed lines refer to the results of the unrestrained editing methods. The solid lines refer to the results of the editing methods coupled with the proposed PRUNE. Statistical significance tests were performed to demonstrate that the improvement in PRUNE compared to baseline was statistically significant (t-test with  $p$ -value  $< 0.05$ ).

Comparing these editing methods, we also observed a notable drop in efficacy and generalization after hundreds of edits with ROME and MEMIT, whereas these values decreased significantly after only 15 edits with MEND. This indicates that in sequential editing scenarios, the MEND method struggled to successfully integrate new knowledge into LLMs after several edits.

**Unrelated facts were perturbed as the number of edits increased.** The locality metric served as an indicator of perturbation for unrelated facts. It became evident that for each editing method, the locality decreased significantly. Additionally, an observation emerged: when the locality of the edited model was low, the performance of downstream tasks was also low. This observation underscores that perturbations of irrelevant knowledge compromise the general abilities of the edited model.

**PRUNE can effectively maintain the editing performance.** This is shown by the solid lines of Figure 4 and could be analyzed from two aspects. On the one hand, when the number of edits was small, the editing performance of each editing method coupled with PRUNE was about the same as the unrestrained method. On the other hand, it significantly mitigated the forgetting of editing facts and the perturbation of irrelevant facts when the number of edits was large during the sequential editing. Specifically, when the number of edits reached 100, the editing performance of MEMIT was very low. But when coupled with PRUNE, its performance remained relatively stable. These observations further validate our motivation in Section 4, demonstrating that the information in the edit update matrices is coordinated, and that performing too many edits can easily result in overfitting. Therefore, applying a certain degree of restraint to edit perturbations can help preserve the model’s general abilities while maintaining the editing knowledge.

### 5.6 EDITING WITH CONCEPTUAL KNOWLEDGE

Section 5.4 and 5.5 analyzed the results on factual knowledge. This section conducted some experiments with ROME on conceptual knowledge using the ConceptEdit dataset (Wang et al., 2024) to make a more comprehensive evaluation. For editing performance, in addition to the three basic metrics, this dataset also designed a new metric “Instance Change” to measure whether the instances under the concept changed accordingly when the definition of the concept was changed.

As shown in Table 2, the performance trends of editing and downstream tasks were similar to those observed with the factual datasets. But there are several key differences: (1) When the number of edits was the same, the editing performance of conceptual knowledge was lower than that of factual knowledge. (2) Both editing performance and general abilities deteriorated more quickly than factual knowledge. For example, even if the number of edits was 100, the editing performance and downstream task performance of ROME were very low, while it was still relatively high when editing factual knowledge. (3) The low “Instance Change” indicated that when the definition of a concept was altered, the instances contained in the original concept were still recognized by the model as belonging to that concept. This shows that this editing method primarily modifies the definition without successfully altering the relationship between concepts and instances, which is not reasonable. These findings indicate that conceptual knowledge is more abstract and more difficult to edit than factual knowledge, highlighting the need to explore editing methods for different types of knowledge.

Table 2: Evaluation results (%) of LLaMA-2 (7B) edited by ROME on the ConceptEdit dataset.

Mode		General Abilities				Editing Performance			
Method	Edits	Reasoning	Summa	Open-QA	NLI	Efficacy	General	Locality	Instance
ROME	20	75.13	11	6.50	24.7	49.15	52.58	35.68	25
	50	20.67	4.90	1.50	0.7	55.42	49.45	19.94	12
	100	12.29	4.7	0.77	0	28.25	30.18	5.68	10
	200	0	4.62	0	0	10.14	8.65	5.31	-8.99
ROME+PRUNE	20	89.38	14.34	23.37	63.54	75.66	58.35	71.7	25
	50	85.15	14.06	25.29	50.52	56.51	45.55	73.16	8
	100	90.78	13.75	21.46	53.17	46.22	42.26	64.06	20
	200	72.9	10.55	22.22	46.15	35.82	34.95	46.65	32

## 5.7 ANALYSIS ON THE FORGETTING OF EDITING FACTS

Section 3 conducted analysis to elucidate the reasons behind the degradation in general abilities with an increasing number of edits. Subsequent experiments quantitatively demonstrated the effectiveness of PRUNE. Here, we delve into qualitative analysis to explain why editing facts are forgotten and how PRUNE can mitigate this forgetting.

Initially, given a set of editing facts  $\mathcal{E} = \{e_1, e_2, \dots\}$ , where  $|\mathcal{E}| = 200$ . ROME was employed for analysis, and the original matrix was defined as  $W$ . During sequential editing, ROME computed key-value pairs  $(k_j^e, v_j^e)$  of the last subject token to generate  $\Delta W_j$  for each edit  $e_j$  to incorporate new facts, satisfying the equation:  $W_j \cdot k_j^e = v_j^e$ . However, when evaluating editing performance, the edited model obtained from the last edit was utilized, thus computing values<sup>5</sup>:  $W_{200} \cdot k_j^e = \hat{v}_j^e$ . After adopting PRUNE to ROME, this equation became  $\overline{W}_{200} \cdot k_j^e = \overline{v}_j^e$ . We hypothesized that if  $\hat{v}_j^e$  was similar to  $v_j^e$ , the editing fact  $e_j$  could be maintained.

Denote  $V_{Current} = \{v_j^e\}$ ,  $V_{Editing} = \{\hat{v}_j^e\}$ , and  $V_{Prune} = \{\overline{v}_j^e\}$ . Specifically, these corresponding values of the first 100 edits were used, as they are more prone to be forgotten than the last 100. Principal Component Analysis (PCA) (Gewers et al., 2022) was employed to visualize these values. The first two principal components of each value were calculated and illustrated, as they can represent most of its features (Zheng et al.). As shown in Figure 5, on the one hand, the discrepancy between the principal components of  $V_{Current}$  and  $V_{Editing}$  was markedly large. This indicates that after 200 edits to the model, the values corresponding to the first 100 facts stored in the edited matrix are severely corrupted, leading to significant forgetfulness. On the other hand, after adopting PRUNE, the discrepancy between the principal components of  $V_{Current}$  and  $V_{Prune}$  was small. This demonstrates that PRUNE effectively maintains the values and mitigates the forgetting of editing facts.

## 6 CONCLUSION AND LIMITATION

In this paper, a theoretical analysis is firstly conducted to elucidate that the bottleneck of the general abilities during sequential editing lies in the escalating value of the condition number. Subsequently, a plug-and-play framework called PRUNE is proposed to apply restraints to preserve general abilities and maintain new editing knowledge simultaneously. Comprehensive experiments on various editing methods and LLMs demonstrate the effectiveness of this method. We aspire that our analysis and method will catalyze future research on continual model editing.

**Limitation** Firstly, this paper focuses on editing a single fact at a time in sequential model editing, but some works study updating hundreds of facts simultaneously in batch editing. Therefore, investigating batch-sequential editing could enhance the scalability of model editing. Secondly, it is necessary to explore the performance of larger-size models and more editing methods on more downstream tasks. Additionally, the proposed PRUNE is only applied once after the last edit. But it could also be utilized multiple times during the sequential editing, and the performance will be better this way.

<sup>5</sup>Since ROME only modifies one matrix, the  $k_j^e$  remains the same across these edited models.

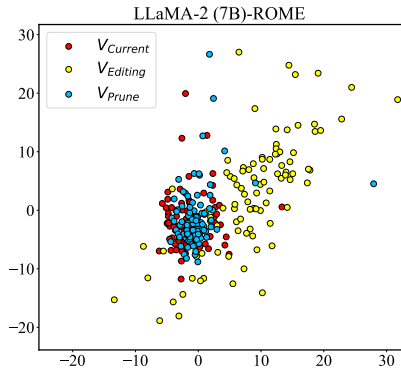


Figure 5: 2-dimensional PCA visualization of first 100 values. The model was edited by ROME with LLaMA-2.

## REFERENCES

- 540  
541  
542 Alfonso M Albano, J Muench, C Schwartz, AI Mees, and PE Rapp. Singular-value decomposition  
543 and the grassberger-procaccia algorithm. *Physical review A*, 38(6):3017, 1988.
- 544 Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh  
545 Sameki, Hanna Wallach, and Kathleen Walker. Fairlearn: A toolkit for assessing and improving  
546 fairness in ai. *Microsoft, Tech. Rep. MSR-TR-2020-32*, 2020.
- 547 Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. In Marie-  
548 Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the*  
549 *2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual*  
550 *Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pp. 6491–6506. Association  
551 for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.522. URL <https://doi.org/10.18653/v1/2021.emnlp-main.522>.
- 552  
553  
554 Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer open-  
555 domain questions. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual*  
556 *Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30*  
557 *- August 4, Volume 1: Long Papers*, pp. 1870–1879. Association for Computational Linguistics,  
558 2017. doi: 10.18653/V1/P17-1171. URL <https://doi.org/10.18653/v1/P17-1171>.
- 559 Shuo Chen, Jindong Gu, Zhen Han, Yunpu Ma, Philip H. S. Torr, and Volker Tresp. Benchmarking  
560 robustness of adaptation methods on pre-trained vision-language models. In Alice Oh,  
561 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.),  
562 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*  
563 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December*  
564 *10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/](http://papers.nips.cc/paper_files/paper/2023/hash/a2a544e43acb8b954dc5846ff0d77ad5-Abstract-Datasets_and_Benchmarks.html)  
565 [hash/a2a544e43acb8b954dc5846ff0d77ad5-Abstract-Datasets\\_and\\_](http://papers.nips.cc/paper_files/paper/2023/hash/a2a544e43acb8b954dc5846ff0d77ad5-Abstract-Datasets_and_Benchmarks.html)  
566 [Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/a2a544e43acb8b954dc5846ff0d77ad5-Abstract-Datasets_and_Benchmarks.html).
- 567 Zhuotong Chen, Zihu Wang, Yifan Yang, Qianxiao Li, and Zheng Zhang. PID control-based self-  
568 healing to improve the robustness of large language models. *CoRR*, abs/2404.00828, 2024. doi: 10.  
569 48550/ARXIV.2404.00828. URL <https://doi.org/10.48550/arXiv.2404.00828>.
- 570  
571 Wenhua Cheng, Weiwei Zhang, Haihao Shen, Yiyang Cai, Xin He, and Kaokao Lv. Optimize weight  
572 rounding via signed gradient descent for the quantization of llms. *CoRR*, abs/2309.05516, 2023.  
573 doi: 10.48550/ARXIV.2309.05516. URL [https://doi.org/10.48550/arXiv.2309.](https://doi.org/10.48550/arXiv.2309.05516)  
574 [05516](https://doi.org/10.48550/arXiv.2309.05516).
- 575 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
576 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
577 Schulman. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168, 2021. URL  
578 <https://arxiv.org/abs/2110.14168>.
- 579  
580 Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects  
581 of knowledge editing in language models. *CoRR*, abs/2307.12976, 2023. doi: 10.48550/ARXIV.  
582 2307.12976. URL <https://doi.org/10.48550/arXiv.2307.12976>.
- 583 Jeff Da, Ronan Le Bras, Ximing Lu, Yejin Choi, and Antoine Bosselut. Analyzing commonsense  
584 emergence in few-shot knowledge models. In Danqi Chen, Jonathan Berant, Andrew McCallum,  
585 and Sameer Singh (eds.), *3rd Conference on Automated Knowledge Base Construction, AKBC*  
586 *2021, Virtual, October 4-8, 2021*, 2021. doi: 10.24432/C5NK5J. URL [https://doi.org/](https://doi.org/10.24432/C5NK5J)  
587 [10.24432/C5NK5J](https://doi.org/10.24432/C5NK5J).
- 588 Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL recognising textual entailment  
589 challenge. In Joaquin Quiñonero Candela, Ido Dagan, Bernardo Magnini, and Florence d’Alché-  
590 Buc (eds.), *Machine Learning Challenges, Evaluating Predictive Uncertainty, Visual Object*  
591 *Classification and Recognizing Textual Entailment, First PASCAL Machine Learning Challenges*  
592 *Workshop, MLCW 2005, Southampton, UK, April 11-13, 2005, Revised Selected Papers*, volume  
593 3944 of *Lecture Notes in Computer Science*, pp. 177–190. Springer, 2005. doi: 10.1007/11736790\_9. URL [https://doi.org/10.1007/11736790\\_9](https://doi.org/10.1007/11736790_9).

- 594 Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons  
595 in pretrained transformers. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.),  
596 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume*  
597 *1: Long Papers)*, ACL 2022, Dublin, Ireland, May 22-27, 2022, pp. 8493–8502. Association for  
598 Computational Linguistics, 2022. doi: 10.18653/v1/2022.acl-long.581. URL <https://doi.org/10.18653/v1/2022.acl-long.581>.
- 600 Jean-Pierre Dedieu. Condition operators, condition numbers, and condition number theorem for the  
601 generalized eigenvalue problem. *Linear algebra and its applications*, 263:1–24, 1997.
- 603 Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh  
604 Ashkboos, Alexander Borzunov, Torsten Hoefer, and Dan Alistarh. Spqr: A sparse-quantized  
605 representation for near-lossless LLM weight compression. *CoRR*, abs/2306.03078, 2023. doi: 10.  
606 48550/ARXIV.2306.03078. URL <https://doi.org/10.48550/arXiv.2306.03078>.
- 607 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts  
608 from diffusion models. *CoRR*, abs/2303.07345, 2023. doi: 10.48550/ARXIV.2303.07345. URL  
609 <https://doi.org/10.48550/arXiv.2303.07345>.
- 611 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-  
612 value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih  
613 (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*,  
614 *EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pp. 5484–  
615 5495. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.446.  
616 URL <https://doi.org/10.18653/v1/2021.emnlp-main.446>.
- 617 Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers build  
618 predictions by promoting concepts in the vocabulary space. In Yoav Goldberg, Zornitsa Kozareva,  
619 and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural*  
620 *Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pp.  
621 30–45. Association for Computational Linguistics, 2022. doi: 10.18653/v1/2022.emnlp-main.3.  
622 URL <https://doi.org/10.18653/v1/2022.emnlp-main.3>.
- 623 Felipe L. Gewers, Gustavo R. Ferreira, Henrique Ferraz de Arruda, Filipi Nascimento Silva, Cesar H.  
624 Comin, Diego R. Amancio, and Luciano da Fontoura Costa. Principal component analysis:  
625 A natural approach to data exploration. *ACM Comput. Surv.*, 54(4):70:1–70:34, 2022. doi:  
626 10.1145/3447755. URL <https://doi.org/10.1145/3447755>.
- 627 Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. SAMSum corpus: A human-  
628 annotated dialogue dataset for abstractive summarization. In Lu Wang, Jackie Chi Kit Cheung,  
629 Giuseppe Carenini, and Fei Liu (eds.), *Proceedings of the 2nd Workshop on New Frontiers in*  
630 *Summarization*, pp. 70–79, Hong Kong, China, November 2019. Association for Computational  
631 Linguistics. doi: 10.18653/v1/D19-5409. URL <https://aclanthology.org/D19-5409>.
- 633 Zhuocheng Gong, Jiahao Liu, Jingang Wang, Xunliang Cai, Dongyan Zhao, and Rui Yan.  
634 What makes quantization for large language model hard? an empirical study from the lens  
635 of perturbation. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (eds.),  
636 *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference*  
637 *on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on*  
638 *Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver,*  
639 *Canada*, pp. 18082–18089. AAAI Press, 2024. doi: 10.1609/AAAI.V38I16.29765. URL  
640 <https://doi.org/10.1609/aaai.v38i16.29765>.
- 641 Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun  
642 Peng. Model editing harms general abilities of large language models: Regularization to the rescue.  
643 In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*,  
644 2024. URL <https://doi.org/10.48550/arXiv.2401.04700>.
- 645 Akshat Gupta and Gopala Anumanchipalli. Rebuilding ROME : Resolving model collapse during  
646 sequential model editing. *CoRR*, abs/2403.07175, 2024. doi: 10.48550/ARXIV.2403.07175. URL  
647 <https://doi.org/10.48550/arXiv.2403.07175>.

- 648 Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. Model editing at scale leads to gradual and  
649 catastrophic forgetting. *CoRR*, abs/2401.07453, 2024a. doi: 10.48550/ARXIV.2401.07453. URL  
650 <https://doi.org/10.48550/arXiv.2401.07453>.  
651
- 652 Akshat Gupta, Dev Sajnani, and Gopala Anumanchipalli. A unified framework for model editing.  
653 *CoRR*, abs/2403.14236, 2024b. doi: 10.48550/ARXIV.2403.14236. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2403.14236)  
654 [10.48550/arXiv.2403.14236](https://doi.org/10.48550/arXiv.2403.14236).
- 655 Anshita Gupta, Debanjan Mondal, Akshay Krishna Sheshadri, Wenlong Zhao, Xiang Lorraine  
656 Li, Sarah Wiegrefe, and Niket Tandon. Editing commonsense knowledge in GPT. *CoRR*,  
657 abs/2305.14956, 2023. doi: 10.48550/ARXIV.2305.14956. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2305.14956)  
658 [48550/arXiv.2305.14956](https://doi.org/10.48550/arXiv.2305.14956).
- 659 Frederik Harder, Matthias Bauer, and Mijung Park. Interpretable and differentially private predictions.  
660 In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second*  
661 *Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI*  
662 *Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA,*  
663 *February 7-12, 2020*, pp. 4083–4090. AAAI Press, 2020. doi: 10.1609/AAAI.V34I04.5827. URL  
664 <https://doi.org/10.1609/aaai.v34i04.5827>.
- 665 Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi.  
666 Aging with GRACE: lifelong model editing with discrete key-value adaptors. In Alice  
667 Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine  
668 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*  
669 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*  
670 *16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/95b6e2ff961580e03c0a662a63a71812-Abstract-Conference.html)  
671 [95b6e2ff961580e03c0a662a63a71812-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/95b6e2ff961580e03c0a662a63a71812-Abstract-Conference.html).
- 672 Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. Does localization inform editing?  
673 surprising differences in causality-based localization vs. knowledge editing in language models.  
674 *CoRR*, abs/2301.04213, 2023. doi: 10.48550/ARXIV.2301.04213. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2301.04213)  
675 [10.48550/arXiv.2301.04213](https://doi.org/10.48550/arXiv.2301.04213).
- 676 Chenhui Hu, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Wilke: Wise-layer knowledge editor  
677 for lifelong knowledge editing. *CoRR*, abs/2402.10987, 2024. doi: 10.48550/ARXIV.2402.10987.  
678 URL <https://doi.org/10.48550/arXiv.2402.10987>.  
679
- 680 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong  
681 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large  
682 language models: Principles, taxonomy, challenges, and open questions. *CoRR*, abs/2311.05232,  
683 2023. doi: 10.48550/ARXIV.2311.05232. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2311.05232)  
684 [2311.05232](https://doi.org/10.48550/arXiv.2311.05232).
- 685 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea  
686 Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Comput.*  
687 *Surv.*, 55(12):248:1–248:38, 2023. doi: 10.1145/3571730. URL [https://doi.org/10.](https://doi.org/10.1145/3571730)  
688 [1145/3571730](https://doi.org/10.1145/3571730).
- 689 Shuoran Jiang, Qingcai Chen, Youcheng Pan, Yang Xiang, Yukang Lin, Xiangping Wu, Chuanyi  
690 Liu, and Xiaobao Song. Zo-adamu optimizer: Adapting perturbation by the momentum and  
691 uncertainty in zeroth-order optimization. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam  
692 Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-*  
693 *Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth*  
694 *Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024,*  
695 *Vancouver, Canada*, pp. 18363–18371. AAAI Press, 2024. doi: 10.1609/AAAI.V38I16.29796.  
696 URL <https://doi.org/10.1609/aaai.v38i16.29796>.
- 697 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris  
698 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion  
699 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav  
700 Petrov. Natural questions: a benchmark for question answering research. *Trans. Assoc. Comput.*  
701 *Linguistics*, 7:452–466, 2019. doi: 10.1162/TACL\_A\_00276. URL [https://doi.org/10.](https://doi.org/10.1162/tacl_a_00276)  
[1162/tacl\\_a\\_00276](https://doi.org/10.1162/tacl_a_00276).

- 702 Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised  
703 open domain question answering. In Anna Korhonen, David R. Traum, and Lluís Màrquez  
704 (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics,*  
705 *ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 6086–6096.  
706 Association for Computational Linguistics, 2019. doi: 10.18653/V1/P19-1612. URL <https://doi.org/10.18653/v1/p19-1612>.  
707
- 708 Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. Zero-shot relation extraction via  
709 reading comprehension. In Roger Levy and Lucia Specia (eds.), *Proceedings of the 21st Conference*  
710 *on Computational Natural Language Learning (CoNLL 2017), Vancouver, Canada, August 3-4,*  
711 *2017*, pp. 333–342. Association for Computational Linguistics, 2017. doi: 10.18653/v1/K17-1034.  
712 URL <https://doi.org/10.18653/v1/K17-1034>.  
713
- 714 Hui-Jia Li, Lin Wang, Yan Zhang, and Matjaž Perc. Optimization of identifiability for efficient  
715 community detection. *New Journal of Physics*, 22(6):063035, 2020.
- 716 Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. Unveiling the  
717 pitfalls of knowledge editing for large language models. *CoRR*, abs/2310.02129, 2023. doi: 10.  
718 48550/ARXIV.2310.02129. URL <https://doi.org/10.48550/arXiv.2310.02129>.  
719
- 720 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization*  
721 *Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.  
722 URL <https://aclanthology.org/W04-1013>.
- 723 Zihao Lin, Mohammad Beigi, Hongxuan Li, Yufan Zhou, Yuxiang Zhang, Qifan Wang, Wenpeng Yin,  
724 and Lifu Huang. Navigating the dual facets: A comprehensive evaluation of sequential memory  
725 editing in large language models. *CoRR*, abs/2402.11122, 2024. doi: 10.48550/ARXIV.2402.11122.  
726 URL <https://doi.org/10.48550/arXiv.2402.11122>.
- 727 Zhi-Quan Luo and Paul Tseng. Perturbation analysis of a condition number for linear systems. *SIAM*  
728 *Journal on Matrix Analysis and Applications*, 15(2):636–660, 1994.  
729
- 730 Jun-Yu Ma, Jia-Chen Gu, Zhen-Hua Ling, Quan Liu, and Cong Liu. Untying the reversal curse via  
731 bidirectional language model editing. *CoRR*, abs/2310.10322, 2023. doi: 10.48550/ARXIV.2310.  
732 10322. URL <https://doi.org/10.48550/arXiv.2310.10322>.
- 733 Jun-Yu Ma, Jia-Chen Gu, Ningyu Zhang, and Zhen-Hua Ling. Neighboring perturbations of  
734 knowledge editing on large language models. *CoRR*, abs/2401.17623, 2024. doi: 10.48550/  
735 ARXIV.2401.17623. URL <https://doi.org/10.48550/arXiv.2401.17623>.  
736
- 737 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
738 associations in GPT. In *NeurIPS*, 2022. URL <https://arxiv.org/abs/2202.05262>.
- 739 Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. Mass-  
740 editing memory in a transformer. In *The Eleventh International Conference on Learning*  
741 *Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL  
742 <https://openreview.net/pdf?id=MkbcAHIYgyS>.
- 743 Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. Fast model  
744 editing at scale. In *The Tenth International Conference on Learning Representations, ICLR 2022,*  
745 *Virtual Event, April 25-29, 2022*. OpenReview.net, 2022a. URL [https://openreview.net/](https://openreview.net/forum?id=0DcZxeWfOPT)  
746 [forum?id=0DcZxeWfOPT](https://openreview.net/forum?id=0DcZxeWfOPT).  
747
- 748 Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. Memory-  
749 based model editing at scale. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári,  
750 Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022,*  
751 *17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning*  
752 *Research*, pp. 15817–15831. PMLR, 2022b. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v162/mitchell122a.html)  
753 [v162/mitchell122a.html](https://proceedings.mlr.press/v162/mitchell122a.html).
- 754 Judea Pearl. Direct and indirect effects. In Jack S. Breese and Daphne Koller (eds.), *UAI*  
755 *'01: Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence, University*  
*of Washington, Seattle, Washington, USA, August 2-5, 2001*, pp. 411–420. Morgan Kaufmann,

2001. URL [https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article\\_id=126&proceeding\\_id=17](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=126&proceeding_id=17).
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *CoRR*, abs/2302.12813, 2023. doi: 10.48550/arXiv.2302.12813. URL <https://doi.org/10.48550/arXiv.2302.12813>.
- Bin Qin, Fu-Lai Chung, and Shitong Wang. KAT: A knowledge adversarial training method for zero-order takagi-sugeno-kang fuzzy classifiers. *IEEE Trans. Cybern.*, 52(7):6857–6871, 2022. doi: 10.1109/TCYB.2020.3034792. URL <https://doi.org/10.1109/TCYB.2020.3034792>.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Chandan Singh, Jeevana Priya Inala, Michel Galley, Rich Caruana, and Jianfeng Gao. Rethinking interpretability in the era of large language models. *CoRR*, abs/2402.01761, 2024. doi: 10.48550/ARXIV.2402.01761. URL <https://doi.org/10.48550/arXiv.2402.01761>.
- Russell A Smith. The condition numbers of the matrix eigenvalue problem. *Numerische Mathematik*, 10:232–240, 1967.
- Gilbert W Stewart and Ji-guang Sun. Matrix perturbation theory. (*No Title*), 1990.
- Ji-guang Sun. Condition number and backward error for the generalized singular value decomposition. *SIAM Journal on Matrix Analysis and Applications*, 22(2):323–341, 2000.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, et al. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023. doi: 10.48550/arXiv.2307.09288. URL <https://doi.org/10.48550/arXiv.2307.09288>.
- Richard J Vaccaro. A second-order perturbation expansion for the svd. *SIAM Journal on Matrix Analysis and Applications*, 15(2):661–671, 1994.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart M. Shieber. Investigating gender bias in language models using causal mediation analysis. In Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Michael E Wall, Andreas Rechtsteiner, and Luis M Rocha. Singular value decomposition and principal component analysis. In *A practical approach to microarray data analysis*, pp. 91–109. Springer, 2003.
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Guozhou Zheng, and Huajun Chen. Easyedit: An easy-to-use knowledge editing framework for large language models. *CoRR*, abs/2308.07269, 2023. doi: 10.48550/arXiv.2308.07269. URL <https://doi.org/10.48550/arXiv.2308.07269>.

- 810 Xiaohan Wang, Shengyu Mao, Ningyu Zhang, Shumin Deng, Yunzhi Yao, Yue Shen, Lei Liang,  
811 Jinjie Gu, and Huajun Chen. Editing conceptual knowledge for large language models. *CoRR*,  
812 abs/2403.06259, 2024. doi: 10.48550/ARXIV.2403.06259. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2403.06259)  
813 [48550/arXiv.2403.06259](https://doi.org/10.48550/arXiv.2403.06259).
- 814 Per-Åke Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numerical*  
815 *Mathematics*, 12:99–111, 1972.
- 816
- 817 Suhang Wu, Minlong Peng, Yue Chen, Jinsong Su, and Mingming Sun. Eva-kellm: A new benchmark  
818 for evaluating knowledge editing of llms. *CoRR*, abs/2308.09954, 2023. doi: 10.48550/ARXIV.  
819 2308.09954. URL <https://doi.org/10.48550/arXiv.2308.09954>.
- 820
- 821 Wanli Yang, Fei Sun, Xinyu Ma, Xun Liu, Dawei Yin, and Xueqi Cheng. The butterfly effect of model  
822 editing: Few edits can trigger large language models collapse. In Lun-Wei Ku, Andre Martins,  
823 and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics, ACL 2024,*  
824 *Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 5419–5437. Association for  
825 Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL.322. URL <https://doi.org/10.18653/v1/2024.findings-acl.322>.
- 826
- 827 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen,  
828 and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. In  
829 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on*  
830 *Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-*  
831 *10, 2023*, pp. 10222–10240. Association for Computational Linguistics, 2023. URL [https://](https://aclanthology.org/2023.emnlp-main.632)  
832 [aclanthology.org/2023.emnlp-main.632](https://aclanthology.org/2023.emnlp-main.632).
- 833
- 834 Lang Yu, Qin Chen, Jie Zhou, and Liang He. MELO: enhancing model editing with neuron-  
835 indexed dynamic lora. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (eds.),  
836 *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on*  
837 *Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational*  
838 *Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pp.  
839 19449–19457. AAAI Press, 2024. doi: 10.1609/AAAI.V38I17.29916. URL [https://doi.](https://doi.org/10.1609/aaai.v38i17.29916)  
[org/10.1609/aaai.v38i17.29916](https://doi.org/10.1609/aaai.v38i17.29916).
- 840
- 841 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi,  
842 Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu,  
843 Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and  
844 Huajun Chen. A comprehensive study of knowledge editing for large language models. *CoRR*,  
845 abs/2401.01286, 2024. doi: 10.48550/ARXIV.2401.01286. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2401.01286)  
[48550/arXiv.2401.01286](https://doi.org/10.48550/arXiv.2401.01286).
- 846
- 847 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo  
848 Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and  
849 Shuming Shi. Siren’s song in the AI ocean: A survey on hallucination in large language  
850 models. *CoRR*, abs/2309.01219, 2023. doi: 10.48550/arXiv.2309.01219. URL [https://](https://doi.org/10.48550/arXiv.2309.01219)  
[doi.org/10.48550/arXiv.2309.01219](https://doi.org/10.48550/arXiv.2309.01219).
- 851
- 852 Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and  
853 Nanyun Peng. On prompt-driven safeguarding for large language models. In *ICLR 2024 Workshop*  
854 *on Secure and Trustworthy Large Language Models*.
- 855
- 856 Zexuan Zhong, Zhengxuan Wu, Christopher D. Manning, Christopher Potts, and Danqi Chen.  
857 Mquake: Assessing knowledge editing in language models via multi-hop questions. In  
858 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on*  
859 *Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-*  
860 *10, 2023*, pp. 15686–15702. Association for Computational Linguistics, 2023. URL [https://](https://aclanthology.org/2023.emnlp-main.971)  
[aclanthology.org/2023.emnlp-main.971](https://aclanthology.org/2023.emnlp-main.971).
- 861
- 862
- 863



## APPENDIX

## A THEORETICAL ANALYSIS BASED ON PERTURBATION THEORY

Here, we provide a detailed analysis and proof of Section 3.2. We begin by introducing some definitions and then present several preliminary lemmas and theorems. These lemmas and theorems are finally used to prove Theorem 3, which is most relevant to our problem discussed in Section 3.2.

## A.1 DEFINITION

We discuss the problem  $Ax = b$ , where  $\tilde{A}$  is a perturbation of  $A$  given by  $\tilde{A} = A + E$ . We assume  $b$  remains unchanged and  $\tilde{x}$  represents the corresponding change, satisfying  $\tilde{A}\tilde{x} = b$ . Here  $A \in \mathbb{C}^{m \times n}$ ,  $b \in \mathbb{C}^m$ .

It is noteworthy that in the following derivation,  $A^H$  denotes the conjugate transpose of  $A$ ,  $A^\dagger$  represents the generalized inverse of  $A$ , and  $\|*\|$  represents 2-norm (Stewart & Sun, 1990).

To simplify the problem, we apply a rotation. Specifically, let  $V = (V_1 \ V_2)$  be a unitary matrix with  $R(V_1) = R(A^H)$ , and let  $U = (U_1 \ U_2)$  be a unitary matrix with  $R(U_1) = R(A)$ , where  $R$  refers to the rank. Then

$$U^H AV = \begin{pmatrix} U_1^H AV_1 & U_1^H AV_2 \\ U_2^H AV_1 & U_2^H AV_2 \end{pmatrix} = \begin{pmatrix} A_{11} & 0 \\ 0 & 0 \end{pmatrix}, \quad (7)$$

where  $A_{11}$  is square and nonsingular. If we set

$$U^H EV = \begin{pmatrix} U_1^H EV_1 & U_1^H EV_2 \\ U_2^H EV_1 & U_2^H EV_2 \end{pmatrix} = \begin{pmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}, \quad (8)$$

then

$$U^H \tilde{A} V = \begin{pmatrix} A_{11} + E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix} = \begin{pmatrix} \tilde{A}_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}. \quad (9)$$

We will call these transformed, partitioned matrices the **reduced form** of the problem. Many statements about the original problem have revealing analogues in the reduced form.

In this form,  $x$  is replaced by  $V^H x$  and  $b$  is replaced by  $U^H b$ . If  $x$  and  $b$  are partitioned in the forms

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad b = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}, \quad (10)$$

where  $x_1, b_1 \in \mathbb{C}^r$ , then

$$x_1 = A_{11}^{-1} b_1 \quad (11)$$

and

$$x_2 = 0. \quad (12)$$

Moreover, the norm of the residual vector

$$r = b - Ax \quad (13)$$

is given by

$$\|r\| = \|b_2\|. \quad (14)$$

Here, we define the symbol  $\eta$ :

$$\eta = \frac{\|A\| \|x\|}{\|b\|}, \quad (15)$$

and for any  $F \in \mathbb{C}^{k \times r}$  ( $k \geq r$ ) the symbol  $\Psi(F)$ , for the spectral norm:

$$\Psi_2(F) = \frac{\|F\|}{(1 + \|F\|^2)^{1/2}}. \quad (16)$$

918 A.2 PRELIMINARY LEMMAS & THEOREMS  
919

920 After introducing some definitions, we give some preliminary lemmas and theorems, which are used  
921 to prove Theorem 3.

922 **Lemma 1** Let

$$\kappa(A) = \|A\| \|A^{-1}\|$$

923 be the condition number of  $A$ . If  $\tilde{A}$  is nonsingular, then

$$\frac{\|\tilde{A}^{-1} - A^{-1}\|}{\|\tilde{A}^{-1}\|} \leq \kappa(A) \frac{\|E\|}{\|A\|}. \quad (17)$$

929 If in addition

$$\frac{\|E\|}{\|A\|} \kappa(A) < 1,$$

932 then  $\tilde{A}$  is performe nonsingular and

$$\|\tilde{A}^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \kappa(A) \frac{\|E\|}{\|A\|}}. \quad (18)$$

937 Moreover

$$\frac{\|\tilde{A}^{-1} - A^{-1}\|}{\|A^{-1}\|} \leq \frac{\kappa(A) \frac{\|E\|}{\|A\|}}{1 - \kappa(A) \frac{\|E\|}{\|A\|}}. \quad (19)$$

942 **Lemma 2** In the reduced form the matrices  $A$  and  $\tilde{A}$  are acute if and only if  $A_{11}$  is nonsingular and

$$E_{22} = E_{21} \tilde{A}_{11}^{-1} E_{12}. \quad (20)$$

945 In this case, if we set

$$F_{21} = E_{21} \tilde{A}_{11}^{-1} \quad \text{and} \quad F_{12} = \tilde{A}_{11}^{-1} E_{12},$$

948 then

$$\tilde{A} = \begin{pmatrix} I \\ F_{21} \end{pmatrix} \tilde{A}_{11} \begin{pmatrix} I & F_{12} \end{pmatrix}$$

951 and

$$\tilde{A}^\dagger = \begin{pmatrix} I & F_{12} \end{pmatrix}^\dagger \tilde{A}_{11}^{-1} \begin{pmatrix} I \\ F_{21} \end{pmatrix}^\dagger. \quad (21)$$

955 **Lemma 3** The matrix

$$\begin{pmatrix} I \\ F \end{pmatrix}$$

958 satisfies

$$\left\| \begin{pmatrix} I \\ F \end{pmatrix}^\dagger \right\| \leq 1 \quad (22)$$

962 and

$$\left\| \begin{pmatrix} I \\ F \end{pmatrix}^\dagger - \begin{pmatrix} I & 0 \end{pmatrix} \right\| = \Psi_2(F). \quad (23)$$

966 **Theorem 1** Let  $\tilde{A}$  be an acute perturbation of  $A$ , and let

$$\hat{\kappa} = \|A\| \|\tilde{A}_{11}^{-1}\|. \quad (24)$$

970 Then

$$\frac{\|\tilde{A}^\dagger - A^\dagger\|}{\|A^\dagger\|} \leq \hat{\kappa} \frac{\|E_{11}\|}{\|A\|} + \Psi_2 \left( \frac{\hat{\kappa} E_{12}}{\|A\|} \right) + \Psi_2 \left( \frac{\hat{\kappa} E_{21}}{\|A\|} \right). \quad (25)$$

972 *Proof.* Let

$$973 \quad I_{21} = \begin{pmatrix} I \\ 0 \end{pmatrix}, \quad I_{12} = (I \quad 0), \quad (26)$$

$$974 \quad J_{21} = \begin{pmatrix} I \\ F_{21} \end{pmatrix}, \quad J_{12} = (I \quad F_{12}). \quad (27)$$

975  
976  
977  
978  $\tilde{A}^\dagger = J_{12}^\dagger A_{11}^{-1} I_{21}^\dagger$ , hence

$$979 \quad \tilde{A}^\dagger - A^\dagger = (J_{12}^\dagger - I_{12}^\dagger) A_{11}^{-1} I_{21}^\dagger + J_{12}^\dagger A_{11}^{-1} (J_{21}^\dagger - I_{21}^\dagger) + J_{12}^\dagger (\tilde{A}_{11}^{-1} - A_{11}^{-1}) J_{21}^\dagger. \quad (28)$$

980  
981  
982 From Lemma 1 we have the following bound:

$$983 \quad \|J_{12}^\dagger (\tilde{A}_{11}^{-1} - A_{11}^{-1}) J_{21}^\dagger\| \leq \|A_{11}^{-1}\| \hat{\kappa} \frac{\|E_{11}\|}{\|A_{11}\|}. \quad (29)$$

984  
985  
986 By Lemma 3

$$987 \quad \|(J_{12}^\dagger - I_{12}^\dagger) A_{11}^{-1} I_{21}^\dagger\| \leq \|A_{11}^{-1}\| \|J_{12}^\dagger - I_{12}^\dagger\| = \|A_{11}^{-1}\| \Psi_2(F_{12}) \quad (30)$$

$$988 \quad = \|A_{11}^{-1}\| \Psi_2(\tilde{A}_{11}^{-1} E_{12}) \quad (31)$$

$$989 \quad \leq \|A_{11}^{-1}\| \Psi_2\left(\frac{\hat{\kappa} E_{12}}{\|A\|}\right), \quad (32)$$

990  
991  
992 and likewise

$$993 \quad \|J_{12}^\dagger A_{11}^{-1} (J_{21}^\dagger - I_{21}^\dagger)\| \leq \|A_{11}^{-1}\| \Psi_2\left(\frac{\hat{\kappa} E_{21}}{\|A\|}\right) \leq \|A^\dagger\| \Psi_2\left(\frac{\hat{\kappa} E_{21}}{\|A\|}\right). \quad (33)$$

994  
995  
996  
997  
998  
999  $\square$

1000  
1001 **Theorem 2** In Theorem 1, let

$$1002 \quad \kappa = \|A\| \|A^\dagger\|, \quad (34)$$

1003 and suppose that

$$1004 \quad \|A^\dagger\| \|E_{11}\| < 1, \quad (35)$$

1005  
1006 so that

$$1007 \quad \gamma \equiv 1 - \frac{\kappa \|E_{11}\|}{\|A\|} > 0. \quad (36)$$

1008  
1009 Then

$$1010 \quad \|\tilde{A}^\dagger\| \leq \frac{\|A^\dagger\|}{\gamma}, \quad (37)$$

1011  
1012 and

$$1013 \quad \frac{\|\tilde{A}^\dagger - A^\dagger\|}{\|A^\dagger\|} \leq \frac{\kappa \|E_{11}\|}{\gamma \|A\|} + \Psi_2\left(\frac{\kappa E_{21}}{\gamma \|A\|}\right) + \Psi_2\left(\frac{\kappa E_{12}}{\gamma \|A\|}\right). \quad (38)$$

1014  
1015  
1016 *Proof.* From the equation  $\tilde{A}^\dagger = J_{12}^\dagger \tilde{A}_{11}^{-1} J_{21}^\dagger$ , we have

$$1017 \quad \|\tilde{A}^\dagger\| \leq \|J_{12}^\dagger\| \|\tilde{A}_{11}^{-1}\| \|J_{21}^\dagger\| \leq \|\tilde{A}_{11}^{-1}\|. \quad (39)$$

1018  
1019  
1020 By Lemma 1,

$$1021 \quad \|\tilde{A}_{11}^{-1}\| \leq \frac{\|A_{11}^{-1}\|}{\gamma} = \frac{\|A^\dagger\|}{\gamma}, \quad (40)$$

1022 which establishes equation 37. Also  $\hat{\kappa} \leq \frac{\kappa}{\gamma}$ , and the inequality equation 38 follows from equation 25.

1023  
1024  
1025  $\square$

1026 A.3 CORE THEOREM  
1027

1028 Finally, we give the core theorem used in main paper. Some symbols and definitions have been  
1029 claimed in Appendix A.1 and A.2.

1030 **Theorem 3** Let  $x = A^\dagger b$  and  $\tilde{x} = \tilde{A}^\dagger b$ , where  $\tilde{A} = A + E$ , and  $E$  is an acute perturbation of  $A$ .  
1031 Then

$$1032 \frac{\|x - \tilde{x}\|}{\|x\|} \leq \hat{\kappa} \frac{\|E_{11}\|}{\|A\|} + \Psi_2 \left( \frac{\hat{\kappa} E_{12}}{\|A\|} \right) + \hat{\kappa}^2 \frac{\|E_{12}\|}{\|A\|} \left( \eta^{-1} \frac{\|b_2\|}{\|b_1\|} + \frac{\|E_{21}\|}{\|A\|} \right). \quad (41)$$

1033  
1034 *Proof.* By Lemma 2, write

$$1035 \tilde{x} - x = J_{12}^\dagger (\tilde{A}_{11}^{-1} - A_{11}^{-1}) b_1 + (J_{12}^\dagger - I_{12}^\dagger) A_{11}^{-1} b_1 + J_{12}^\dagger \tilde{A}_{11}^{-1} (J_{21}^\dagger - I_{21}^\dagger) b. \quad (42)$$

1036  
1037 Then

$$1038 \|J_{12}^\dagger (\tilde{A}_{11}^{-1} - A_{11}^{-1}) b_1\| \leq \hat{\kappa} \frac{\|E_{11}\|}{\|A\|} \|x\|, \quad (43)$$

1039  
1040 and

$$1041 \|(J_{12}^\dagger - I_{12}^\dagger) A_{11}^{-1} b_1\| \leq \Psi_2 \left( \frac{\hat{\kappa} E_{12}}{\|A\|} \right) \|x\|. \quad (44)$$

1042  
1043 Now

$$1044 J_{12}^\dagger \tilde{A}_{11}^{-1} (J_{21}^\dagger - I_{21}^\dagger) b = J_{12}^\dagger \tilde{A}_{11}^{-1} ((I + F_{21}^H F_{21})^{-1} - I) b_1 + J_{12}^\dagger \tilde{A}_{11}^{-1} (I + F_{21}^H F_{21})^{-1} F_{21}^H b_2. \quad (45)$$

1045  
1046 To bound the first term in equation 45, note that

$$1047 (I + F_{21}^H F_{21})^{-1} - I = -(I + F_{21}^H F_{21})^{-1} F_{21}^H F_{21}.$$

1048  
1049 Hence

$$1050 \begin{aligned} 1051 \|J_{12}^\dagger \tilde{A}_{11}^{-1} ((I + F_{21}^H F_{21}) - I) b_1\| &\leq \|\tilde{A}_{11}^{-1}\| \| (I + F_{21}^H F_{21})^{-1} \| \|F_{21}^H\| \|F_{21} b_1\| \\ 1052 &\leq \|\tilde{A}_{11}^{-1}\| \|E_{21}\| \|\tilde{A}_{11}^{-1} b_1\| \\ 1053 &\leq \|\tilde{A}_{11}^{-1}\| \|E_{21}\|^2 \|x\| \\ 1054 &= \left( \frac{\hat{\kappa} \|E_{21}\|_2}{\|A\|} \right)^2 \|x\|. \end{aligned} \quad (46)$$

1055  
1056 For the second term in equation 45 we have

$$1057 \begin{aligned} 1058 \|J_{12}^\dagger \tilde{A}_{11}^{-1} (I + F_{21}^H F_{21})^{-1} F_{21} b_2\| &\leq \|\tilde{A}_{11}^{-1}\|^2 \|E_{21}\| \|b_2\| \\ 1059 &= \|\tilde{A}_{11}^{-1}\|^2 \|E_{21}\| \frac{\|b_2\|}{\|b_1\|} \eta^{-1} \|x\| \|A\| \\ 1060 &\leq \eta^{-1} \hat{\kappa}^2 \frac{\|E_{21}\| \|b_2\|}{\|A\| \|b_1\|} \|x\|. \end{aligned} \quad (47)$$

1061  
1062 The bound equation 41 follows on combining equation 42–equation 47.  
1063  
1064  
1065  
1066  
1067  $\square$

1068 Readers can refer to this work (Stewart & Sun, 1990) for more details of perturbation analysis.

1069 Returning to our problem, consider  $Wk = v$ , where  $(k, v) \in P$ . Let  $\tilde{W} = W + \Delta W$ , where  $\Delta W$   
1070 is the corresponding perturbation matrix. Assuming  $v$  remains constant, there exists  $\Delta k$  such that  
1071  $\tilde{k} = k + \Delta k$  satisfies  $\tilde{W}\tilde{k} = v$ . And we have  $k = W^\dagger v$  and  $\tilde{k} = \tilde{W}^\dagger v$ . Applying Theorem 3, we  
1072 obtain

$$1073 \frac{\|\Delta k\|}{\|k\|} = \frac{\|k - \tilde{k}\|}{\|k\|} \leq \hat{\kappa} \frac{\|\Delta E_{11}\|}{\|W\|} + \Psi_2 \left( \frac{\hat{\kappa} \Delta E_{12}}{\|W\|} \right) + \hat{\kappa}^2 \frac{\|\Delta E_{12}\|}{\|W\|} \left( \eta^{-1} \frac{\|v_2\|}{\|v_1\|} + \frac{\|\Delta E_{21}\|}{\|W\|} \right), \quad (48)$$

1074 where  $\Delta E_{11}$ ,  $\Delta E_{12}$ ,  $\Delta E_{21}$ , and  $\Delta W$  are directly related, and each term on the right-hand side  
1075 involves  $\hat{\kappa}$ . This means that the relative perturbation of the vector  $k$  is constrained by  $\hat{\kappa}$ . According  
1076 to Theorem 2,  $\hat{\kappa} \leq \frac{\kappa}{\gamma}$ , where  $\kappa = \|W\| \|W^\dagger\|$  is the condition number of  $W$ . This indicates that  $\kappa$  is  
1077 a robust indicator of the impact of  $\Delta W$  on the vector  $k$ .  
1078  
1079

## B EXPERIMENTAL SETUP

### B.1 BASELINE EDITING METHODS

Three popular model editing methods were selected as baselines including:

- **MEND** (Mitchell et al., 2022a)<sup>6</sup>: it learned a hypernetwork to produce weight updates by decomposing the fine-tuning gradients into rank-1 form.
- **ROME** (Meng et al., 2022)<sup>7</sup>: it first localized the factual knowledge at a specific layer in the transformer MLP modules, and then updated the knowledge by directly writing new key-value pairs in the MLP module.
- **MEMIT** (Meng et al., 2023)<sup>8</sup>: it extended ROME to edit a large set of facts and updated a set of MLP layers to update knowledge.

The ability of these methods were assessed based on EasyEdit<sup>9</sup> (Wang et al., 2023), an easy-to-use knowledge editing framework which integrates the released codes and hyperparameters from previous methods.

### B.2 EDITING DATASETS AND EVALUATION METRICS

Table 3 shows the examples of two factual datasets (ZSRE) (Levy et al., 2017) and COUNTERFACT (Meng et al., 2022). Figure 6 shows an example of ConceptEdit dataset, which is cited from Wang et al. (2024). More details can refer to the original paper of these datasets.

Table 3: The editing datasets of both ZSRE and COUNTERFACT.

Datasets	Editing prompt
ZSRE	Which was the record label for New Faces, New Sounds?
COUNTERFACT	In America, the official language is

**An Example of ConceptEdit Dataset**

Concept	publisher	Phrased prompt	When we refer to publisher, we are talking about
Descriptor x	The definition of publisher is	Locality sentence	The definition of curling league is a group of sports teams that compete against each other ...
Descriptor y (intra module)	a group of sports teams or individual athletes that ...	Instances	['Famitsu Bunko', 'BitComposer', ...]
Descriptor y (inter module)	very tall building	Instance prompt	Whether Famitsu Bunko belongs to category publisher?

Figure 6: An example of ConceptEdit dataset

Besides, following previous works (Meng et al., 2022; Mitchell et al., 2022a; Meng et al., 2023), the editing performance metrics for the ZSRE and COUNTERFACT datasets are efficacy, generalization and locality, but there are some computational differences. In the main paper, the metrics of editing performance are used for the ZSRE dataset.

For the COUNTERFACT dataset, here are the details:

**Efficacy** validates whether the edited models could recall the editing fact under editing prompt  $p_i$ . The assessment is based on Efficacy Score (**ES**) representing as:  $\mathbb{E}_i[\mathbb{1}[P_{f_{\theta_n}}(o_i^* | p_i) > P_{f_{\theta_n}}(o_i | p_i)]]$ , where  $\mathbb{1}$  is the indicator function.

<sup>6</sup><https://github.com/eric-mitchell/mend>

<sup>7</sup><https://github.com/kmeng01/rome>

<sup>8</sup><https://github.com/kmeng01/memit>

<sup>9</sup><https://github.com/zjunlp/EasyEdit>

**Generalization** verifies whether the edited models could recall the editing fact under the paraphrase prompts  $\mathcal{P}_i^G$  via Generalization Score (**GS**):  $\mathbb{E}_i [\mathbb{E}_{p \in \mathcal{P}_i^G} [\mathbb{1}[P_{f_{\theta_n}}(o_i^* | p) > P_{f_{\theta_n}}(o_i | p)]]]$ .

**Locality** verifies whether the output of the edited models for inputs out of editing scope remains unchanged under the locality prompts  $\mathcal{P}_i^L$  via Locality Score (**LS**):  $\mathbb{E}_i [\mathbb{E}_{p_l \in \mathcal{P}_i^L} [\mathbb{1}[P_{f_{\theta_n}}(o_l | p_l) > P_{f_{\theta_n}}(o_i^* | p_l)]]]$ , where  $o_l$  was the original answer of  $p_l$ .

### B.3 HYPERPARAMETERS OF PRUNE

When conducting experiments, for different editing methods, LLMs and editing datasets, the hyperparameter  $\alpha$  in function  $F$  of PRUNE is different. Table 4 shows the details of this hyperparameter.  $e$  is the base of the natural logarithm.

Table 4: The hyperparameters  $\alpha$  for PRUNE.

Datasets	Models	ROME	MEMIT	MEND
COUNTERFACT	GPT-2 XL	1.2	1.2	1.2
	LLaMA-2	1.2	$e$	1.2
	LLaMA-3	1.5	$e$	-
ZsRE	LLaMA-2	1.2	$e$	$e$

### B.4 TASK PROMPTS

The prompts for each downstream task were illustrated in Table 5.

Table 5: The prompts to LLMs for evaluating their zero-shot performance on these general tasks.

Reasoning:

Q: {QUESTION} A: Let’s think step by step. {HINT} Therefore, the answer (arabic numerals) is:

NLI:

{SENTENCE1} entails the {SENTENCE2}. True or False? answer:

Open-domain QA:

Refer to the passage below and answer the following question. Passage: {DOCUMENT} Question: {QUESTION}

Summarization:

{DIALOGUE} TL;DR:

### B.5 EXPERIMENTS COMPUTE RESOURCES

We used NVIDIA A800 80GB GPU for experiments. For LLaMA-2 (7B) and LLaMA-3 (8B), it occupies about 40+GB memory and costs about 3 hours for each editing method to run 200 edits and then to test downstream tasks. For GPT-2 XL (1.5B), it needs 10+GB and costs about 1.5 hours for each editing method to run 200 edits and then to test downstream tasks.

## C EXPERIMENTAL RESULTS

### C.1 RESULTS OF GENERAL ABILITIES

Figure 7, 8 and 9 show the downstream task performance of edited models with GPT-2 XL, LLaMA-2 (7B) and LLaMA-3 (8B) on COUNTERFACT dataset. Due to limitations of computing resources, experiments were conducted using only LLaMA-2 (7B) on the ZsRE dataset. We will supplement experiments with other LLMs in the future.

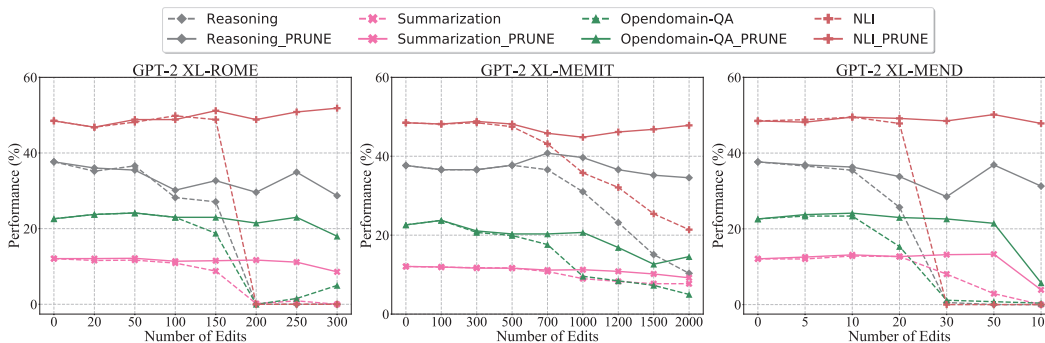


Figure 7: The downstream task performance (%) of models edited by three editing methods with GPT-2 XL on the COUNTERFACT dataset.

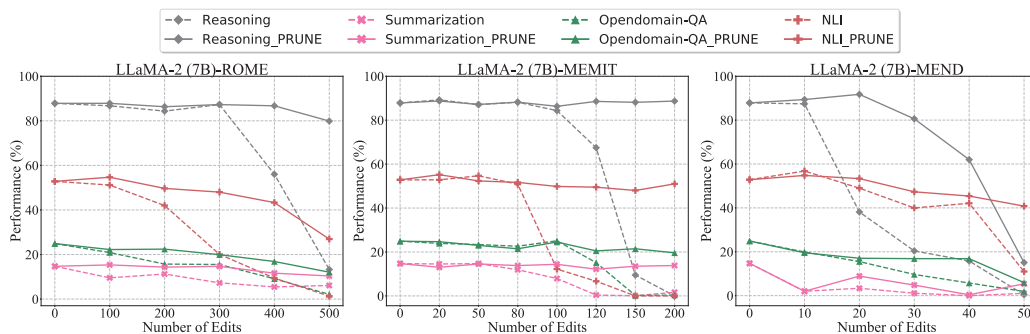


Figure 8: The downstream task performance (%) of models edited by three editing methods with LLaMA-2 (7B) on the COUNTERFACT dataset.

## C.2 RESULTS OF EDITING PERFORMANCE

Figure 10, 11 and 12 shows the editing performance of edited models with GPT-2 XL, LLaMA-2 (7B) and LLaMA-3 (8B) on COUNTERFACT dataset.

## C.3 RESULTS OF ANOTHER FUNCTION FOR PRUNE

In the main paper, log function is used in  $F$  in PRUNE to restrain  $\hat{\sigma}_i$ . Here we use the linear function, which could be represented as:  $F(\hat{\sigma}_i) = \frac{1}{\beta} * \hat{\sigma}_i + \frac{\beta-1}{\beta} * \max\{\sigma_i\}$ . Here  $\beta > 1$  was a hyperparameter and was set as 2 in this section. Figure 13 and 14 respectively show some downstream task performance and editing performance with linear function on COUNTERFACT dataset.

Compared with Figure 7 and 10, we observed that although the linear function in PRUNE played a role in preserving general abilities and maintaining editing performance, its effectiveness was noticeably inferior to that of the log function when the number of edits was large.

## C.4 CONDITION NUMBER WITH PRUNE

Figure 15 shows after coupling with PRUNE, the condition number of MEMIT is significantly restrained.

## C.5 THE CORRELATION BETWEEN CONDITION NUMBER AND GENERAL ABILITIES

Figure 16 simultaneously shows the condition number and general abilities of three editing methods without PRUNE in the sequential editing process. From these experiments, we observed that a dramatic increase in the condition number is often accompanied by a rapid decline in general abilities.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

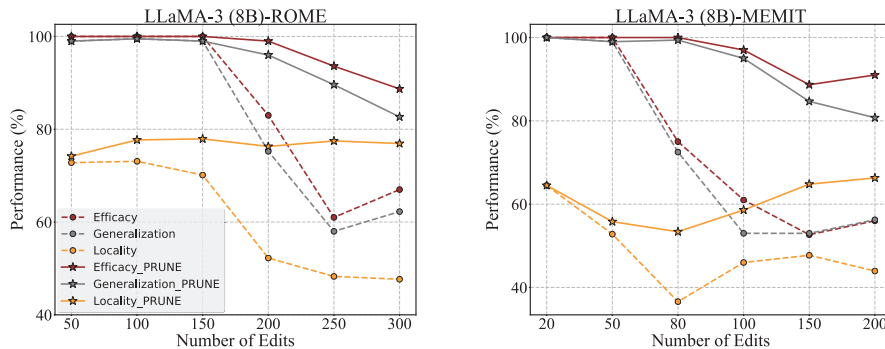


Figure 9: The downstream task performance (%) of models edited by two editing methods with LLaMA-3 (8B) on the COUNTERFACT dataset. Since the code framework EasyEdit used in this paper does not currently support MEND editing on LLaMA-3, there are no results of MEND here.

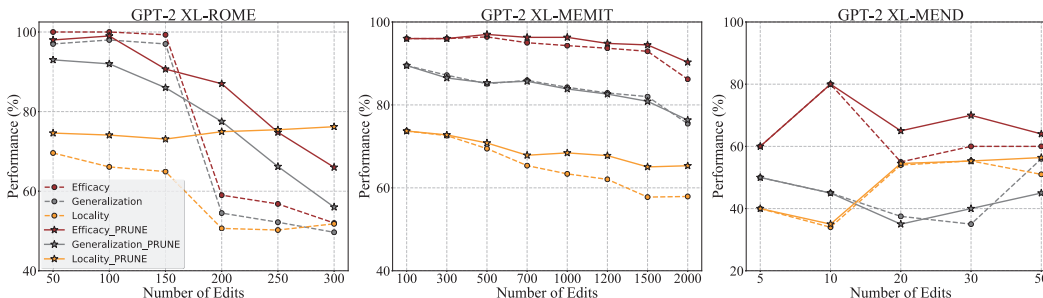


Figure 10: The editing performance (%) of three editing methods with GPT-2 XL on COUNTERFACT dataset.

## D BROADER IMPACTS

This work offers significant advancements in the field of model editing for LLMs. By addressing the challenge of preserving general abilities while performing sequential edits, PRUNE facilitates continual learning and adaptability in LLMs. This can lead to several positive impacts, such as:

**Enhanced Adaptability.** It enables LLMs to update their knowledge base quickly and accurately without extensive retraining. This adaptability is crucial in dynamic environments where up-to-date information is vital, such as real-time translation services, personalized learning systems, and interactive virtual assistants.

**Resource Efficiency.** By mitigating the need for full retraining, PRUNE significantly reduces computational resources and energy consumption. This aligns with sustainable AI and makes it more feasible to deploy LLMs in resource-constrained settings.

**Improved Performance in Specialized Tasks.** PRUNE’s ability to perform targeted edits without compromising overall model performance can enhance LLMs’ effectiveness in specialized domains, such as medical diagnostics, legal analysis, and technical support, where precise and updated knowledge is essential.

While this work offers many benefits, there are potential negative societal impacts that must be considered:

**Misuse for Malicious Purposes.** The capability to edit LLMs efficiently could be exploited to inject harmful or biased information into models, thereby spreading disinformation or propaganda. This risk is particularly concerning in applications involving social media and news dissemination where LLMs might generate or amplify misleading content.



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

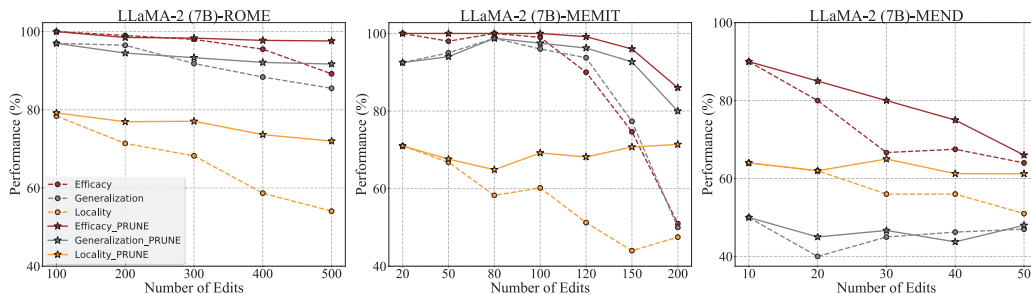


Figure 11: The editing performance (%) of three editing methods with LLaMA-2 (7B) on the COUNTERFACT dataset.

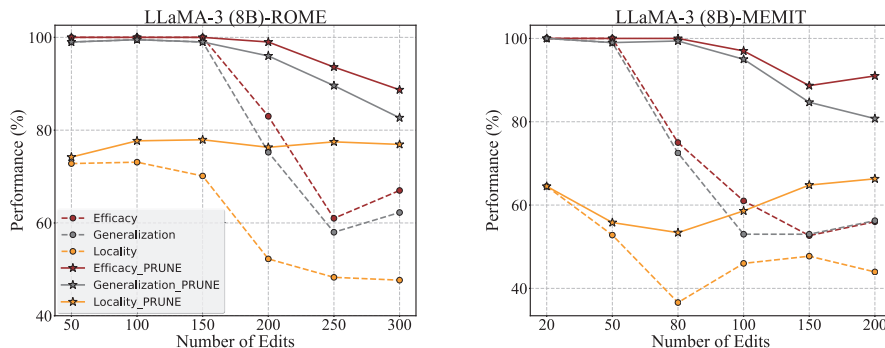


Figure 12: The editing performance (%) of three editing methods with LLaMA-3 (8B) on the COUNTERFACT dataset.

**Fairness.** Unintended biases could be introduced during the editing process, potentially exacerbating existing biases in LLMs. This could lead to unfair treatment or misrepresentation of specific groups, especially if the editing is not conducted with proper oversight and consideration of ethical implications.

**Privacy Concerns.** The ability to update models quickly might also pose privacy risks, as models could be edited to include sensitive or personal information. Ensuring that editing processes do not compromise individual privacy is critical, particularly in applications involving personal data.

To mitigate these potential negative impacts, several strategies could be implemented:

**Gated Release and Monitoring.** Limiting access to the framework through gated releases and monitoring its usage can help prevent misuse.

**Bias and Fairness Audits.** Conducting regular audits to assess and address biases in the model editing process can help ensure that edits do not unfairly impact any specific group. Developing guidelines for ethical editing practices is also essential.

**Privacy Protection Measures.** Establishing clear protocols for handling sensitive data during the editing process can help protect privacy. Anonymization and encryption techniques should be employed to safeguard personal information.

By considering both the positive and negative impacts and implementing appropriate mitigation strategies, this work can contribute to the responsible and ethical advancement of model editing technologies.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361

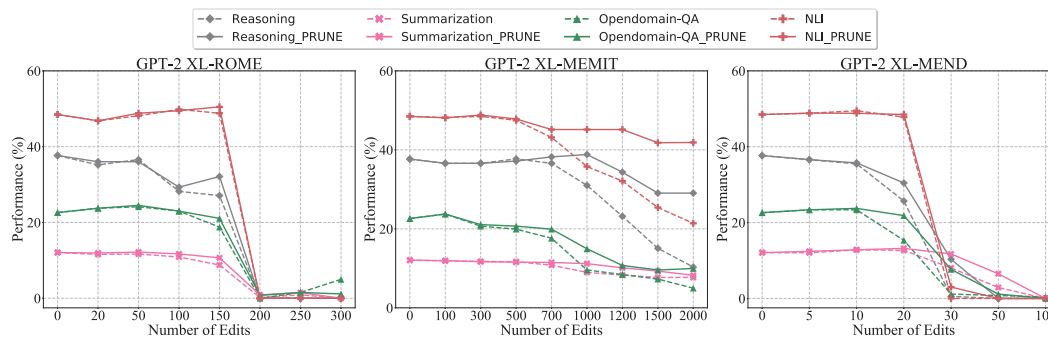


Figure 13: The downstream task performance (%) of models edited by three editing methods with GPT-2 XL on the COUNTERFACT dataset. Here the linear function was used in PRUNE.

1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374

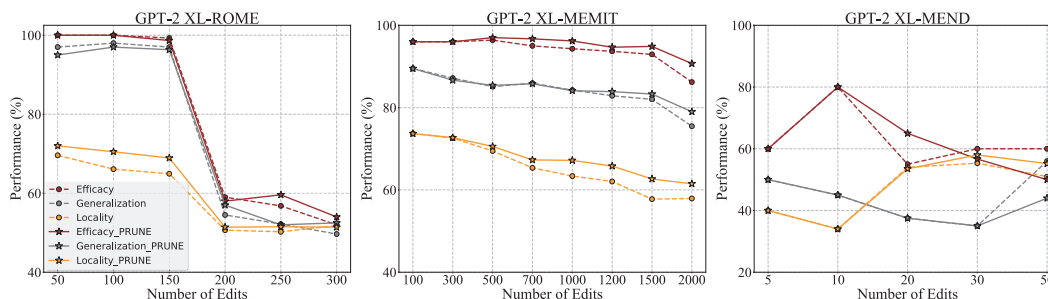


Figure 14: The editing performance (%) of editing methods with GPT-2 XL on the COUNTERFACT dataset. Here the linear function was used in PRUNE.

1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387

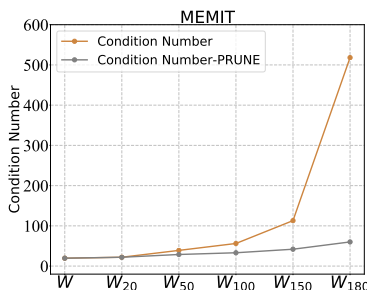


Figure 15: The condition number of MEMIT with LLaMA-2 (7B) on the COUNTERFACT dataset. “-PRUNE” refers to the condition number of MEMIT coupled with the proposed PRUNE.

1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401

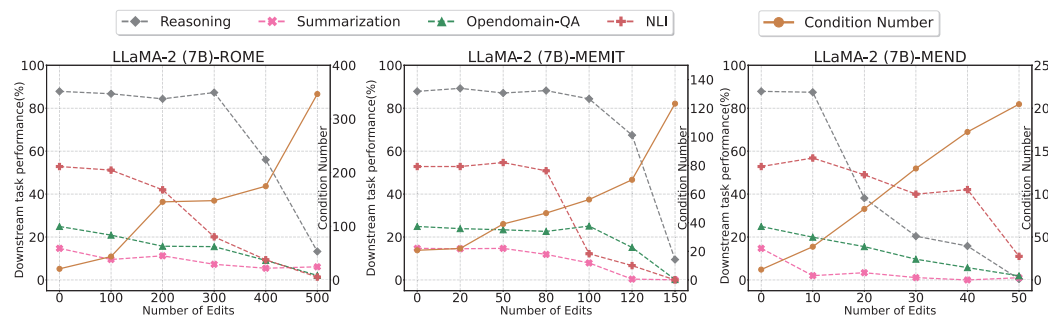


Figure 16: The condition number and downstream task performance of three editing methods with LLaMA-2 (7B) on the COUNTERFACT dataset. Since MEMIT and MEND have multiple parameters to be edited, we randomly selected one of them to calculate the condition number.

1402  
1403