

Neuron-Level Sequential Editing for Large Language Models

Anonymous ACL submission

Abstract

This work explores sequential model editing in large language models (LLMs), a critical task that involves modifying internal knowledge within LLMs continuously through multi-round editing, each incorporating updates or corrections to adjust the model’s outputs without the need for costly retraining. Existing model editing methods, especially those that alter model parameters, typically focus on single-round editing and often face significant challenges in sequential model editing—most notably issues of model forgetting and failure. To address these challenges, we introduce a new model editing method, namely Neuron-level Sequential Editing (NSE), tailored for supporting sequential model editing. Specifically, we optimize the target layer’s hidden states using the model’s original weights to prevent model failure. Furthermore, we iteratively select neurons in multiple layers for editing based on their activation values to mitigate model forgetting. Our empirical experiments demonstrate that NSE significantly outperforms current modifying parameters model editing methods, marking a substantial advancement in the field of sequential model editing. Our code is released on <https://anonymous.4open.science/r/NSE-0A8D/>.

1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities in storing extensive factual knowledge during pre-training and recalling this information during inference (Brown et al., 2020; Petroni et al., 2019; Roberts et al., 2020). However, as real-world knowledge continuously evolves, the information within these models can become outdated or incorrect (Cao et al., 2021; Mitchell et al., 2022a; Fang et al., 2024). Retraining LLMs to incorporate new information is often prohibitively costly (Mitchell et al., 2022b; Meng et al., 2022). Consequently, recent years have witnessed a surge

in *model editing* methods focusing on modifying specific knowledge without the complete retraining process. Specifically, they first identify the crucial layers for the target knowledge by calculating their casual effect on output. Then, by updating the weights of these layers, they manipulate the hidden states of these layers to modify the final output, allowing LLMs to seamlessly adapt to dynamic real-world information (Meng et al., 2023; Hartvigsen et al., 2023).

Although current direct model editing methods are effective for single-round modifications, real-world applications require a continuous learning process where models must retain previous edits during subsequent modifications (Yao et al., 2023). This has led to the concept of *sequential model editing*, which requires the performance of multiple consecutive edits on models. However, current direct model editing methods pose significant risks in this context (Meng et al., 2022, 2023). The primary risk is *model forgetting*, where cumulative changes in parameters from consecutive edits cause the model to forget previously modified knowledge, thus degrading overall performance (Gupta et al., 2024a). For example, as illustrated in Figure 1 (a), after editing the model with new knowledge about “The cat”, the LLM forgets previously edited knowledge about “The latest Olympic”.

Furthermore, the second risk is *model failure*, where excessive edits impair the model’s ability to generate coherent text. Worse still, this impairment can potentially lead to model collapse, characterized by producing irrelevant, repetitive, or non-sensical text, as illustrated in Figure 1 (a). In view of this, recent research, such as memory-based methods (Mitchell et al., 2022b; Hartvigsen et al., 2023), has attempted to address these challenges by preserving the LLM parameters after each edit. However, the increasing storage requirements as the number of edits grows significantly limit the practicality of these methods.

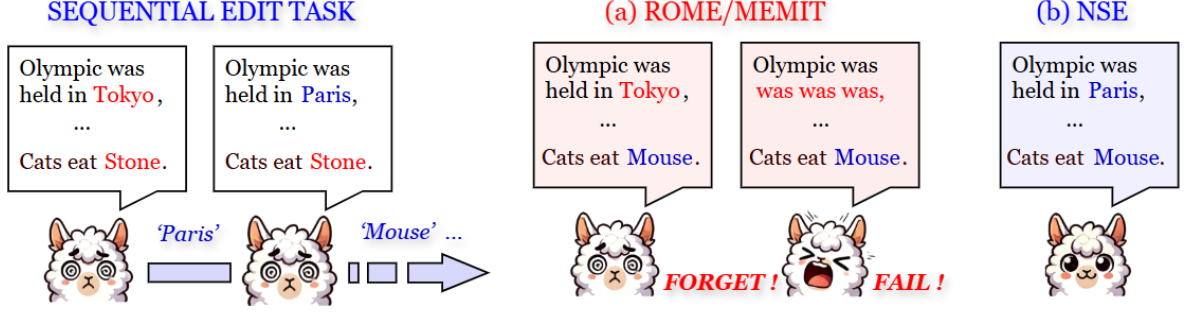


Figure 1: Example of sequential editing. (a) shows model forgetting and model failure issues in sequential editing using ROME/MEMIT, while (b) shows the accurate editing capabilities of our method without such issues.

To tackle these challenges, we introduce a new model editing method, termed **Neuron-level Sequential Editing (NSE)**. Specifically, to address the model failure, NSE uses weights rewinding for value computation by preserving the model’s original weights as a significant reference when manipulating the hidden states of the crucial layers. This process can effectively mitigate the impairment of previous knowledge accumulated over various edits. Furthermore, to address model forgetting, NSE selectively collects “influential neurons” to update weights by sorting neuron activation within the crucial layers, rather than updating all weights in the critical layers as in previous work (Meng et al., 2022, 2023). This selective modification maximizes the protection of model functionality from degradation. Additionally, for large-scale LLMs containing numerous neurons, an iterative multi-layer editing is introduced to streamline the neuron selection process, enabling NSE to achieve massive knowledge updates effectively in a single editing.

Through theoretical analysis and extensive experiments conducted on GPT2-XL (1.5B) (Radford et al., 2019), GPT-J (6B) (Wang and Komatsuzaki, 2021), and Llama3 (8B), we validate the effectiveness and efficiency of our NSE. Compared to current model editing methods (e.g., MEND (Mitchell et al., 2022a), ROME (Meng et al., 2022), MEMIT (Meng et al., 2023) and GRACE (Hartvigsen et al., 2023)), NSE shows substantial improvements *with respect to* five commonly used metrics such as specificity and consistency.

2 Preliminary

2.1 Autoregressive Language Model

An autoregressive language model predicts the next token in a sequence based on the tokens that have come before it. Given a L layer transformer model

and an input sequence $x = (x_0, x_1, \dots, x_T)$, the model aims to predict the next token in the sequence. The probability of the next token x_{t+1} is given by:

$$\mathbb{P}(x_{t+1} | x_0, x_1, \dots, x_t) = \text{Softmax}(\mathbf{W}_e \mathbf{h}_t^N), \quad (1)$$

where \mathbf{W}_e represents embedding matrix, and \mathbf{h}_t^N represents the final hidden state at the topmost layer N . The hidden state \mathbf{h}_t^l at layer l is calculated as:

$$\begin{aligned} \mathbf{h}_t^l(x) &= \mathbf{h}_t^{l-1}(x) + \mathbf{a}_t^l(x) + \mathbf{v}_t^l(x), \\ \mathbf{a}_t^l &= \text{attn}^l(\mathbf{h}_0^{l-1}, \mathbf{h}_1^{l-1}, \dots, \mathbf{h}_t^{l-1}), \\ \mathbf{v}_t^l &= \mathbf{W}_{\text{out}}^l \sigma(\mathbf{W}_{\text{in}}^l \gamma(\mathbf{h}_t^{l-1} + \mathbf{a}_t^l)), \end{aligned} \quad (2)$$

where \mathbf{a}_t^l represents the output of the attention block and \mathbf{v}_t^l represents the output of the FFN layer. \mathbf{W}_{in}^l and $\mathbf{W}_{\text{out}}^l$ are weight matrices, σ is a non-linear activation function, and γ represents the layer norm.

2.2 Previous Model Editing Method

Sequential model editing refines a pre-trained model f_{θ_0} through multiple updates, each incorporating new facts (s, r, o) , such as $s =$ “The latest Olympic”, $r =$ “was held in”, $o =$ “Paris”. After each edit, the updated model f_{θ_t} is optimized to produce the correct outputs for the current edit set $\mathbb{D}_{\text{edit}_t}$ while maintaining the precision of previous tasks, ensuring both adaptation to new information and preservation of previous performance.

Following current work (Meng et al., 2022, 2023), we treat the weights of the Transformer’s (Vaswani et al., 2017) FFN layer as a linear associative memory. That is, linear operations within the FFN layer can be viewed as key-value storage for information retrieval (Kohonen, 1972; Anderson, 1972).

Given the weights \mathbf{W}_{in}^l of the l -th FFN layer when prompted by (s_i, r_i) , we identify the activation output of the last subject token S as the key \mathbf{k}_i^l . In the following, this key \mathbf{k}_i^l is processed through the output weights $\mathbf{W}_{\text{out}}^l$, producing the value \mathbf{v}_i^l . In the context of sequential model editing, we start with an initial set of key-value associations for knowledge facts stored in the l -th FFN layer, denoted respectively by $\mathbf{K}_0 = \{\mathbf{k}_i\}_{i=1}^n$ and $\mathbf{V}_0 = \{\mathbf{v}_i\}_{i=1}^n$. In the following, our objective is to introduce new key value associations m , denoted as $\mathbf{K}_1 = \{\mathbf{k}_i\}_{i=n+1}^{n+m}$ and $\mathbf{V}_1 = \{\mathbf{v}_i\}_{i=n+1}^{n+m}$, while retaining all existing associations unchanged. Drawing on prior work (Meng et al., 2023), the optimization objective is as follows:

$$\Delta^* = \arg \min_{\Delta} (\|(\mathbf{W} + \Delta)\mathbf{K}_1 - \mathbf{V}_1\|^2 + \|(\mathbf{W} + \Delta)\mathbf{K}_0 - \mathbf{V}_0\|^2), \quad (3)$$

where \mathbf{W} represents the weights of \mathbf{W}_{out} in the target FFN layer, Δ denotes the update to \mathbf{W} , and \mathbf{V}_1 can be directly trained using the fine-tuning loss predicted by the model through backpropagation. Since \mathbf{K}_0 and \mathbf{V}_0 represent the knowledge retained in LLM, we can express this relationship as $\mathbf{W}\mathbf{K}_0 = \mathbf{V}_0$. Therefore, we can obtain the closed-form solution of Eqn. 3 using the least squares method (Lang, 2012):

$$\Delta^* = \mathbf{R}\mathbf{K}_1^T (\mathbf{K}_0\mathbf{K}_0^T + \mathbf{K}_1\mathbf{K}_1^T)^{-1}, \quad (4)$$

where $\mathbf{R} = \mathbf{V}_1 - \mathbf{W}\mathbf{K}_1$. Furthermore, MEMIT focuses on editing a specific set of layers denoted as $\mathcal{R} = \{l_0 - |\mathcal{R}| + 1, \dots, l_0\}$. The required the weight update Δ^l for layer $l \in \mathcal{R}$ is expressed as:

$$\Delta^l = \mathbf{R}^l \mathbf{K}_1^{lT} (\mathbf{K}_0^l \mathbf{K}_0^{lT} + \mathbf{K}_1^l \mathbf{K}_1^{lT})^{-1}, \quad (5)$$

where $\mathbf{R}^l = \frac{\mathbf{R}^{l_0}}{l_0 - l + 1}$. These modifications are implemented sequentially, starting from the lower layers and progressing to the upper layers.

3 Methodology

In this section, we introduce NSE, a method adapted for sequential model editing, as illustrated in Figure 2. Initially, we identify the values computation method in Section 3.1. Subsequently, in Section 3.2, we detail the editing method that selectively filters certain neurons for the corresponding parameter updates. Finally, in Section 3.3, we introduce the approach of iterative multi-layer editing.

3.1 Weights Rewinding for Value Computation

First, our primary goal is to find a hidden vector that encodes the new association (s_i, r_i, o_i) and replaces the value \mathbf{v}_i^l in the l -th layer as described in Section 2.2. In practical implementation, as shown in Figure 2 (a), we optimize δ_i through gradient descent by maximizing the probability that the model outputs o_i to compute $\mathbf{z}_i = \mathbf{h}_i^l + \delta_i$, where \mathbf{h}_i^l denotes the hidden state of the LLM in layer l . And the value \mathbf{v}_i^l can be calculated as $\mathbf{v}_i^l + \delta_i$. In the process of sequential editing, we observed that using the updated model parameters f_{θ_t} to calculate \mathbf{z}_i after each edit leads to significant model degradation over multiple edit rounds. This indicates that the cumulative parameter updates from each editing round can lead to a change in value computation. In contrast, using the original model parameters f_{θ_0} to compute \mathbf{z}_i effectively prevents this issue. Hence, we propose a weight rewinding method for the value computation, which is based on the initial model weights f_{θ_0} to ensure that \mathbf{z}_i is computed using f_{θ_0} for each edit. The optimization objective is as follows:

$$\delta_i^* = \arg \min_{\delta_i} -\log \mathbb{P}_{f_{\theta_0}(\mathbf{h}_i^l + \delta_i)} [o_i | (s_i, r_i)], \quad (6)$$

$$\mathbf{z}_i = \mathbf{h}_i^l + \delta_i^*,$$

where $f_{\theta_0}(\mathbf{h}_i^l + \delta_i)$ represents the original model with \mathbf{h}_i^l updated to $\mathbf{h}_i^l + \delta_i$. Subsequently, the value \mathbf{v}_i^l can be updated to $\mathbf{v}_i^l + \delta_i^*$, which will be used to support subsequent model edits. Note that when calculating \mathbf{z}_i using the original model parameters f_{θ_0} , it is sufficient to save only the weight matrix \mathbf{W}_{out} that needs to be updated, rather than storing the entire model parameters and we use the original weights only during the computation of \mathbf{z}_i .

3.2 Neuron-level Weights Updating

In Section 3.1, we primarily discussed the method to calculate \mathbf{z}_i to replace the value in the target layer. In this section, we will specifically elaborate on how to utilize the computed value for weight updates at the neuron level, as shown in Figure 2 (b).

Based on previous work, it is established that neurons in FFN contain abundant information (Dai et al., 2022; Wang et al., 2022; Schwettmann et al., 2023; Pan et al., 2023; Wu et al., 2023). Hence, we selectively optimize a subset of neurons rather than

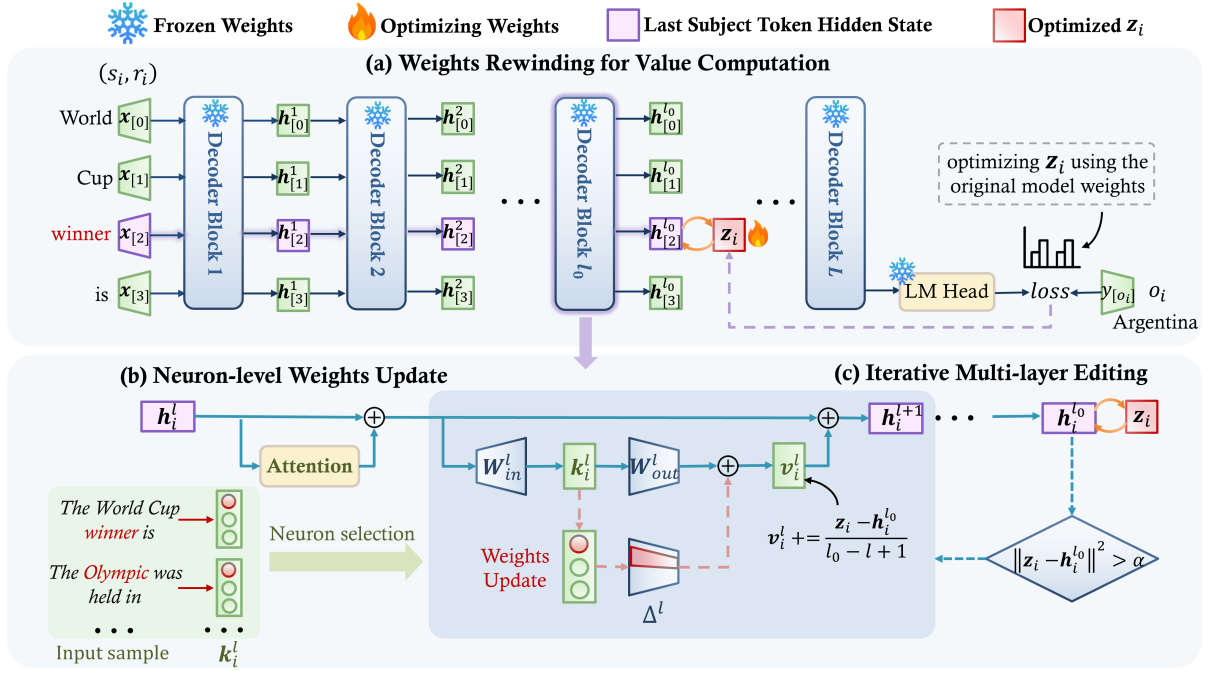


Figure 2: Overview of sequential model editing with NSE. (a) describes the process of weights rewinding for value computation. (b) illustrates the neuron selection and neuron-level weights update. (c) shows the process of iterative multi-layer editing.

alter the entire weight matrix for each edit. Specifically, for a given knowledge fact (s_i, r_i, o_i) , we use the activation values k_i of the neurons and compute scores $Q_i = |k_i|$. Neurons are ranked based on these scores, and a subset is chosen such that the cumulative score of selected neurons surpasses a predetermined percentage of the total score:

$$\begin{aligned} \mathcal{I} = \arg \min_{\mathcal{I} \subseteq \{1, \dots, N\}} |\mathcal{I}| \\ \text{s.t.} \quad \sum_{j \in \mathcal{I}} Q_{ij} \geq p \times \sum_{j=1}^N Q_{ij}, \end{aligned} \quad (7)$$

where Q_{ij} represents the score of the j -th neuron and \mathcal{I} is the set of indices for selected neurons.

Given the need for batch editing, which entails editing multiple knowledge facts simultaneously, each possibly corresponding to different neuron sets, we compute the sum of neuron scores corresponding to all samples in a batch to derive a new neuron score, which is then used for selecting neurons.

Hereafter, we introduce how to update \mathbf{W} (i.e., \mathbf{W}_{out}) by the selected neurons set \mathcal{I} . Let $\hat{\mathbf{W}}$ and $\hat{\Delta}$ denote the submatrices of \mathbf{W} and Δ , respectively, selected according to the indices of the set \mathcal{I} . Our objective is to modify only a subset of neuronal parameters by altering specific weight rows in \mathbf{W} to achieve the optimization described in Eqn. 3.

And we can transform Eqn. 3 accordingly:

$$\begin{aligned} \hat{\Delta}^* = \arg \min_{\hat{\Delta}} (\|(\hat{\mathbf{W}} + \hat{\Delta})\hat{\mathbf{K}}_1 - \mathbf{V}_1\|^2 \\ + \|(\hat{\mathbf{W}} + \hat{\Delta})\hat{\mathbf{K}}_0 - \mathbf{V}_0\|^2), \end{aligned} \quad (8)$$

where $\hat{\mathbf{K}}_0$ and $\hat{\mathbf{K}}_1$ are the submatrices of \mathbf{K}_0 and \mathbf{K}_1 formed by selecting elements indexed by \mathcal{I} respectively. Based on the method of minimal squared error, we solve Eqn. 8 as follows:

$$\hat{\Delta}^* = \hat{\mathbf{R}}\hat{\mathbf{K}}_1^T\hat{\mathbf{C}}^{-1}, \quad (9)$$

where $\hat{\mathbf{R}} = \mathbf{V}_1 - \hat{\mathbf{W}}\hat{\mathbf{K}}_1$ and $\hat{\mathbf{C}} = \hat{\mathbf{K}}_0\hat{\mathbf{K}}_0^T + \hat{\mathbf{K}}_1\hat{\mathbf{K}}_1^T$, and following current work (Meng et al., 2023), $\mathbf{K}_0\mathbf{K}_0^T$ is estimated by $\lambda \mathbb{E}[\mathbf{k}\mathbf{k}^T]$, where λ is a hyperparameter that balances the weights between new knowledge and preserved knowledge. The submatrix $\hat{\mathbf{K}}_0\hat{\mathbf{K}}_0^T$ is then obtained by selecting the rows and columns indexed by \mathcal{I} from $\mathbf{K}_0\mathbf{K}_0^T$. Additionally, during the continuous editing process, the newly edited knowledge from each round will become the old knowledge for the next round. Therefore, after each round, we will add the newly edited knowledge into $\mathbf{K}_0\mathbf{K}_0^T$.

3.3 Iterative multi-layer Editing

As described in Section 2.2, current work (Meng et al., 2023) propagates edits through layers by

computing the value v_i^l as $v_i^l = \frac{\delta_i}{l_0 - l + 1}$ ($l \in \mathcal{R}$) (Meng et al., 2023; Gupta et al., 2024b). Here, $\delta_i = z_i - h_i^{l_0}$ represents the residual difference. The fundamental purpose is that, as each layer is updated, the hidden state $h_i^{l_0}$ progressively approaches the target z_i , thus diminishing the residual δ_i . Detailed analyses can be found in the Appendix E. However, due to errors in the fitting process, some knowledge proves difficult to edit, resulting in $v_i^{l_0}$ not adequately approximating z_i , and consequently leading to editing failures.

Therefore, we propose iterative multi-layer editing to refine the multi-layer editing approach in MEMIT by iteratively selecting neurons to edit multiple layers, as depicted in Figure 2 (c). Specifically, considering that some knowledge is difficult to edit such that the corresponding value $v_i^{l_0}$ cannot sufficiently approximate the optimized target value z_i , we employ iterative multi-layer editing. After each round of multi-layer editing, we filter the knowledge samples in the current batch based on $\|z_i - h_i^{l_0}\|^2$. If $\|z_i - h_i^{l_0}\|^2 < \alpha$, the knowledge sample is considered edited successfully. In contrast, if $\|z_i - h_i^{l_0}\|^2 > \alpha$, the sample is deemed not yet successfully edited. The hyperparameter α , which is set differently for various LLMs, determines the threshold for editing success. These unedited samples are then filtered to form a new batch for further multi-layer editing, repeating this process until all knowledge samples in the batch meet $\|z_i - h_i^{l_0}\|^2 < \alpha$ or the iteration limit is reached. More details of the NSE algorithm are provided in Appendix A.

4 Experiments

We conducted experiments to demonstrate the effectiveness of our model editing method. The experiments aim to address the following research questions:

- **RQ1:** How does NSE perform on sequential editing tasks compared to existing methods?
- **RQ2:** What is the impact of adjusting the batch size of edits on the performance of NSE?
- **RQ3:** Can the LLM, after undergoing NSE editing, retain its original general capabilities, and how does it perform on general capability tests?

4.1 Experimental Settings

Datasets & Evaluation Metrics. To evaluate the effectiveness of our method, we use two datasets: Counterfact (Meng et al., 2022) and ZsRE (Levy

et al., 2017). For the Counterfact dataset, we employ five evaluation metrics as defined in previous work (Meng et al., 2022, 2023): **Efficacy** (efficiency success), **Generalization** (paraphrase success), **Specificity** (neighborhood success), **Fluency** (generation entropy), and **Consistency** (reference score). For the ZsRE dataset, we use three evaluation metrics also defined in previous work (Mitchell et al., 2022a; Meng et al., 2022, 2023): **Efficacy**, **Generalization**, and **Specificity**. For more details, see Appendix C.

Models & Baselines. Our comparative analysis evaluates the performance of various editing methods in three LLMs, GPT2-XL (1.5B) (Radford et al., 2019), GPT-J (6B) (Wang and Komatsuzaki, 2021) and Llama3 (8B). For baseline comparisons, we primarily select model editing methods that modify the model’s parameters, including fine-tuning the specific layer (FT-L) (Zhu et al., 2020), MEND (Mitchell et al., 2022a), ROME (Meng et al., 2022) and MEMIT (Meng et al., 2023). Furthermore, we incorporated memory-based editing methods, SERAC (Mitchell et al., 2022b) and GRACE (Hartvigsen et al., 2023), as a baseline. More details are provided in Appendix B.

4.2 Performance Comparison (RQ1)

In this subsection, we provide a comprehensive comparison of NSE with existing methods on the sequential model editing task, using GPT2-XL, GPT-J, and Llama3 models. The experiments are conducted with a total of 2000 edited samples and an editing batch size of 100 (batch size refers to the number of samples edited simultaneously in each editing round during the sequential editing process), evaluated on the Counterfact and ZsRE datasets. The results of all the evaluation methods, on all datasets, are presented in Table 1. Furthermore, we test the methods on the edited samples after each edit round on the GPT2-XL model using Counterfact. We presented the result on ROME, MEMIT and NSE in Figure 3. We also provide a case study that compares the effects of text generation of different editing methods, with results available in the Appendix F. According to these, we can find that:

- **Observation 1:** NSE outperforms other baseline methods in almost all critical metrics in both datasets and models in the sequential editing task. Specifically, compared to baseline methods for parameter modification, NSE shows significant improvements in all metrics. In par-

Table 1: Comparison of NSE with existing methods on the sequential model editing task. *Eff.*, *Gen.*, *Spe.*, *Flu.* and *Consis.* denote Efficacy, Generalization, Specificity, Fluency and Consistency, respectively.

Method	Model	Counterfact					ZsRE		
		Eff.↑	Gen.↑	Spe.↑	Flu.↑	Consis.↑	Eff.↑	Gen.↑	Spe.↑
Pre-edited		7.85 \pm 0.26	10.58 \pm 0.26	89.48 \pm 0.18	635.23 \pm 0.11	24.14 \pm 0.08	36.99 \pm 0.30	36.34 \pm 0.30	31.89 \pm 0.22
FT-L	Llama3	83.33 \pm 0.37	67.79 \pm 0.40	46.63 \pm 0.37	233.72 \pm 0.22	8.77 \pm 0.05	30.48 \pm 0.26	30.22 \pm 0.32	15.49 \pm 0.17
MEND		63.24 \pm 0.31	61.17 \pm 0.36	45.37 \pm 0.38	372.16 \pm 0.80	4.21 \pm 0.05	0.91 \pm 0.05	1.09 \pm 0.05	0.53 \pm 0.02
ROME		64.40 \pm 0.47	61.42 \pm 0.42	49.44 \pm 0.38	449.06 \pm 0.26	3.31 \pm 0.02	2.01 \pm 0.07	1.80 \pm 0.07	0.69 \pm 0.03
MEMIT		65.65 \pm 0.47	64.65 \pm 0.42	51.56 \pm 0.38	437.43 \pm 1.67	6.58 \pm 0.11	34.62 \pm 0.36	31.28 \pm 0.34	18.49 \pm 0.19
SERAC		71.21 \pm 0.56	61.05 \pm 0.39	66.90 \pm 0.21	615.72 \pm 0.34	20.77 \pm 0.13	67.75 \pm 0.24	33.96 \pm 0.35	22.17 \pm 0.15
GRACE		90.72 \pm 0.13	0.09 \pm 0.01	87.23 \pm 0.21	632.43 \pm 0.63	23.79 \pm 0.23	74.58\pm0.31	1.03 \pm 0.06	31.86 \pm 0.12
NSE		96.14\pm0.19	78.42\pm0.35	87.66\pm0.19	632.72\pm0.12	30.20\pm0.10	62.29 \pm 0.35	47.13\pm0.31	32.32\pm0.22
Pre-edited		22.23 \pm 0.73	24.34 \pm 0.62	78.53 \pm 0.33	626.64 \pm 0.31	31.88 \pm 0.20	22.19 \pm 0.24	31.30 \pm 0.27	24.15 \pm 0.32
FT-L	GPT2-XL	63.55 \pm 0.48	42.20 \pm 0.41	57.06 \pm 0.30	519.35 \pm 0.27	10.56 \pm 0.05	37.11 \pm 0.39	33.30 \pm 0.37	10.36 \pm 0.17
MEND		50.80 \pm 0.50	50.80 \pm 0.48	49.20 \pm 0.51	407.21 \pm 0.08	1.01 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
ROME		54.60 \pm 0.48	51.18 \pm 0.40	52.68 \pm 0.33	366.13 \pm 1.40	0.72 \pm 0.02	47.50 \pm 0.43	43.56 \pm 0.42	14.27 \pm 0.19
MEMIT		94.70 \pm 0.22	85.82 \pm 0.28	60.50 \pm 0.32	477.26 \pm 0.54	22.72 \pm 0.15	79.17 \pm 0.32	71.44 \pm 0.36	26.12 \pm 0.25
SERAC		82.28 \pm 0.26	58.31 \pm 0.34	68.98 \pm 0.32	615.92 \pm 0.72	28.65 \pm 0.17	92.37 \pm 0.29	38.21 \pm 0.32	25.17 \pm 0.25
GRACE		94.50 \pm 0.24	0.04 \pm 0.01	78.13\pm0.43	622.56 \pm 0.79	31.55 \pm 0.25	82.54 \pm 0.21	0.40 \pm 0.02	24.78 \pm 0.21
NSE		96.80\pm0.20	87.72\pm0.30	72.10 \pm 0.28	622.85\pm0.15	40.04\pm0.11	83.26\pm0.29	75.33\pm0.34	26.14\pm0.25
Pre-edited		16.22 \pm 0.31	18.56 \pm 0.45	83.11 \pm 0.13	621.81 \pm 0.67	29.74 \pm 0.51	26.32 \pm 0.37	25.79 \pm 0.25	27.42 \pm 0.53
FT-L	GPT-J	92.15 \pm 0.27	72.38 \pm 0.38	43.35 \pm 0.37	297.92 \pm 0.77	6.65 \pm 0.10	72.37 \pm 0.29	68.91 \pm 0.32	19.66 \pm 0.23
MEND		46.15 \pm 0.50	46.22 \pm 0.51	53.90 \pm 0.48	242.41 \pm 0.41	3.94 \pm 0.03	0.71 \pm 0.04	0.71 \pm 0.04	0.52 \pm 0.03
ROME		57.50 \pm 0.48	54.20 \pm 0.40	52.05 \pm 0.31	589.28 \pm 0.08	3.22 \pm 0.02	56.42 \pm 0.42	54.65 \pm 0.42	9.86 \pm 0.16
MEMIT		98.55 \pm 0.11	95.50\pm0.16	63.64 \pm 0.31	546.28 \pm 0.88	34.89 \pm 0.15	94.91 \pm 0.16	90.22 \pm 0.23	27.56 \pm 0.27
SERAC		72.25 \pm 0.15	58.18 \pm 0.23	64.06 \pm 0.37	595.35 \pm 0.35	27.35 \pm 0.12	92.17 \pm 0.67	36.57 \pm 0.72	20.67 \pm 0.22
GRACE		95.88 \pm 0.28	0.05 \pm 0.01	82.11\pm0.24	620.21 \pm 0.49	28.53 \pm 0.15	94.33 \pm 0.37	1.59 \pm 0.03	27.63 \pm 0.43
NSE		99.55\pm0.06	91.92 \pm 0.22	78.96 \pm 0.25	620.49\pm0.16	40.24\pm0.12	96.87\pm0.14	91.33\pm0.22	28.66\pm0.25

ticular, on Llama3 (8B) editing, NSE achieves an average enhancement of around 30.33% on multiple metrics. Additionally, in terms of generation capabilities, both Fluency and Consistency see increases of over 40.75%. In contrast, while the GRACE parameter preservation method retains model capabilities to the greatest extent, it exhibits weaker generalization performance.

- **Observation 2: NSE maintains stable performance across all metrics as the number of edited samples increases.** As shown in Figure 3, NSE shows robust performance, remains resilient to model failure, and forgets despite the increase in the number of editing rounds. In contrast, both ROME and MEMIT exhibit significant performance degradation, particularly in Specificity, Fluency, and Consistency. This degradation suggests that as the number of edited sam-

ples increases, ROME and MEMIT struggle to maintain model integrity, severely impairing the model’s generative capabilities and leading to progressive model failure and forgetting.

4.3 Impact of Batch Size (RQ2)

To answer RQ2, we examine the impact of different batch sizes on the performance of NSE compared to MEMIT, in the sequential model editing task, with a total of 2000 edits on the Counterfact dataset. Figure 4 presents four radar charts that depict Llama3 performance (8B) when using batch sizes of 200, 100, 50, and 10, respectively. We also tested the editing effects of GPT2-XL and GPT-J at different batch sizes, as detailed in Appendix G. According to Figure 4, we can find that:

- **Observation 3: NSE maintains superior performance across various batch sizes in the se-**

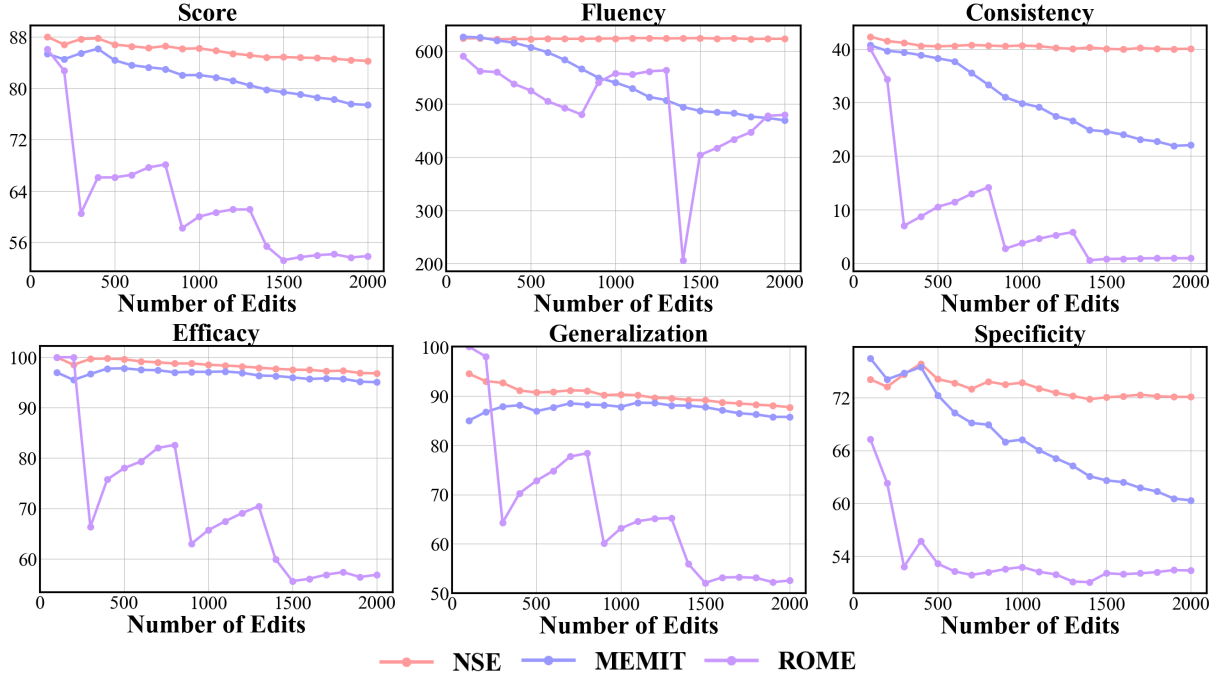


Figure 3: Editing performance of NSE and baselines with varying numbers of edits (batch size 100) in sequential editing, evaluated on the Counterfact dataset. Score is the harmonic mean of Efficacy, Generalization, and Specificity.

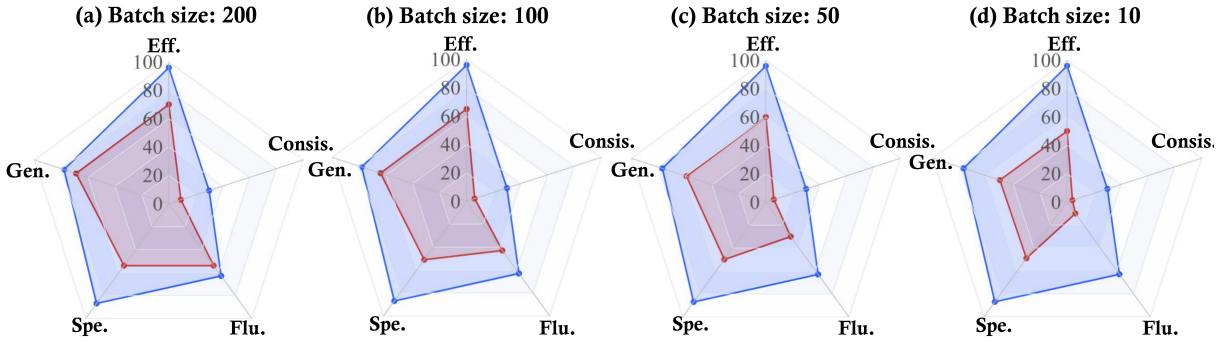


Figure 4: Editing performance of NSE and MEMIT with different batch size, evaluated on Llama3 (8B). The red line and the blue line represent MEMIT and NSE, respectively.

quential editing task. Specifically, as the batch size decreases, resulting in an increase in the number of editing rounds, the performance of MEMIT deteriorates. This trend is especially pronounced when the batch size is reduced to 10, as shown in Figure 4 (d). The radar charts reveal a significant decline in model editing effectiveness, with notable decreases observed in all metrics. In contrast, NSE demonstrates an average improvement in all of these metrics.

4.4 General Ability Test (RQ3)

To assess the effects of model editing on the general capabilities of edited model, we have selected six natural language tasks from the General Language Understanding Evaluation (GLUE) bench-

mark (Wang et al., 2019). Details of the downstream tasks are described in the Appendix C.2.

We perform evaluations on Llama3 (8B) based on sequential editing settings with 3000 edits. The results are shown in Figure 5. Note that as the parameter-preserved editing method differs from ours and does not affect the model by changing it, it will not be compared with it in the general ability test. Here we can make the following observations:

- **Observation 4: NSE consistently maintains the general capabilities of the LLM during sequential editing without incurring model failure.** Specifically, as the number of edited knowledge instances increases, NSE’s performance remains aligned with that of the pre-edited LLM, demonstrating no adverse effects on the LLM’s

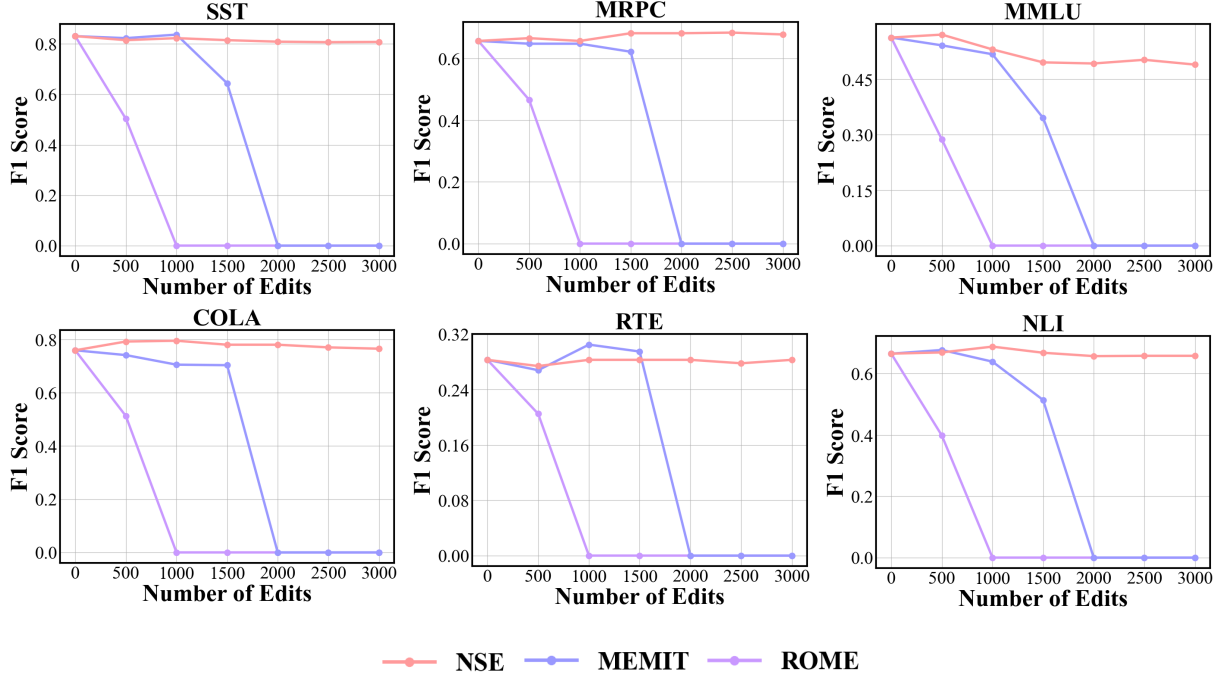


Figure 5: Performance on general tasks of edited models using NSE, ROME and MEMIT, with sequential editing on Llama3 (8B).

inherent general capabilities. In contrast, ROME and MEMIT exhibit a significant decline in general capabilities 0 after editing approximately 1,000 to 2,000 samples, indicating that the model has already experienced degradation.

5 Related work

Current approaches to model editing in LLMs generally fall into two main categories (Yao et al., 2023; Zhang et al., 2024; Bi et al., 2024a):

Preserve Models’ Parameters. Methods that preserve the original model’s parameters generally store edit examples in memory and use them to guide the model’s predictions (Bi et al., 2024b). For example, SERAC (Mitchell et al., 2022b) keeps the original model unchanged and uses a separate counterfactual model for edits. T-Patcher (Huang et al., 2023) introduces an additional neuron for each output error, whereas CaliNet (Dong et al., 2022) incorporates knowledge using a predetermined number of neurons. GRACE (Hartvigsen et al., 2023) implements sequential editing by maintaining a dynamically updated codebook. WISE (Wang et al., 2024) introduces a dual parametric memory scheme that separates pretrained and edited knowledge.

Modify Models’ Parameters. In contrast, methods that modify the parameters of the LLM require updating the model’s internal parameters with each

edit. FT-W fine-tunes the specific layer with regularization constraints (Zhu et al., 2020). KN (Dai et al., 2022) identifies the "knowledge neuron" as the crucial pair in FFN and updates these neurons, which represents the knowledge. KE (Cao et al., 2021) and MEND (Mitchell et al., 2022a) use hypernetwork to predict the weight change of LLM based on meta-learning methods. ROME (Meng et al., 2022) and MEMIT (Meng et al., 2023) support large-scale direct edits by locating and editing knowledge in specific GPT layers. To improve the performance of parameter modification methods in sequential editing, PRUNE (Ma et al., 2024) restricts the maximum singular value of parameter changes. In this paper, we focus mainly on parameter-modification editing methods.

6 Conclusion

In this work, we introduce NSE, a model editing method for sequential model editing tailored to address the significant challenges of model forgetting and model failure. Specifically, we propose weight rewinding for value computation by optimizing the hidden states of the target layer using the model’s original weights, which effectively minimize cumulative changes and maintain model coherence. Additionally, we select influential neurons for different knowledge to update weights in FFN and iteratively edit multi-layer weights.

7 Limitations

Despite the outstanding performance of NSE in sequential editing, our investigation reveals some limitations. Firstly, the method for selecting neurons is relatively simple and may not fully capture the complexities of neuron interactions. Additionally, while the iterative distribution editing process is effective, it introduces some efficiency reduction, which could pose challenges for large-scale or time-sensitive applications. Moving forward, our goal is to explore more effective neuron attribution methods and enhance the efficiency of our editing techniques.

References

- James A Anderson. 1972. A simple neural network generating an interactive memory. *Mathematical biosciences*, 14(3-4):197–220.
- Luisa Bentivogli, Bernardo Magnini, Ido Dagan, Hoa Trang Dang, and Danilo Giampiccolo. 2009. The fifth PASCAL recognizing textual entailment challenge. In *TAC*. NIST.
- Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, and Xueqi Cheng. 2024a. [Is factuality decoding a free lunch for llms? evaluation on knowledge editing benchmark](#). *Preprint*, arXiv:2404.00216.
- Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, Hongcheng Gao, Junfeng Fang, and Xueqi Cheng. 2024b. Struedit: Structured outputs enable the fast and accurate knowledge editing for large language models. *arXiv preprint arXiv:2409.10132*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *NeurIPS*.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. In *EMNLP (1)*, pages 6491–6506. Association for Computational Linguistics.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *ACL (1)*, pages 8493–8502. Association for Computational Linguistics.
- William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases.

- In *IWP@IJCNLP*. Asian Federation of Natural Language Processing.
- Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. Calibrating factual knowledge in pretrained language models. In *EMNLP (Findings)*, pages 5937–5947. Association for Computational Linguistics.
- Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Xiang Wang, Xiangnan He, and Tat seng Chua. 2024. Alphaedit: Null-space constrained knowledge editing for language models. *arXiv preprint arXiv:2410.02355*.
- Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. 2024a. Model editing at scale leads to gradual and catastrophic forgetting. *CoRR*, abs/2401.07453.
- Akshat Gupta, Dev Sajani, and Gopala Anumanchipalli. 2024b. A unified framework for model editing. *CoRR*, abs/2403.14236.
- Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. Aging with GRACE: lifelong model editing with discrete key-value adapters. In *NeurIPS*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In *ICLR*. OpenReview.net.
- Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformer-patcher: One mistake worth one neuron. In *ICLR*. OpenReview.net.
- Teuvo Kohonen. 1972. Correlation matrix memories. *IEEE Trans. Computers*, 21(4):353–359.
- Serge Lang. 2012. *Introduction to linear algebra*. Springer Science & Business Media.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. In *CoNLL*, pages 333–342. Association for Computational Linguistics.
- Jun-Yu Ma, Hong Wang, Hao-Xiang Xu, Zhen-Hua Ling, and Jia-Chen Gu. 2024. Perturbation-restrained sequential model editing. *CoRR*, abs/2405.16821.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In *NeurIPS*.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer. In *ICLR*. OpenReview.net.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. Fast model editing at scale. In *ICLR*. OpenReview.net.

- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022b. Memory-based model editing at scale. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 15817–15831. PMLR.
- Haowen Pan, Yixin Cao, Xiaozhi Wang, and Xun Yang. 2023. Finding and editing multi-modal neurons in pre-trained transformer. *CoRR*, abs/2311.07470.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander H. Miller. 2019. Language models as knowledge bases? In *EMNLP/IJCNLP (1)*, pages 2463–2473. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In *EMNLP (1)*, pages 5418–5426. Association for Computational Linguistics.
- Sarah Schwettmann, Neil Chowdhury, Samuel Klein, David Bau, and Antonio Torralba. 2023. Multimodal neurons in pretrained text-only transformers. In *ICCV (Workshops)*, pages 2854–2859. IEEE.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, pages 1631–1642. ACL.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*, pages 5998–6008.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *ICLR (Poster)*. OpenReview.net.
- Ben Wang and Aran Komatsuzaki. 2021. Gpt-j-6b: A 6 billion parameter autoregressive language model.
- Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-jun Chen. 2024. WISE: rethinking the knowledge memory for lifelong model editing of large language models. *CoRR*, abs/2405.14768.
- Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. Finding skill neurons in pre-trained transformer-based language models. In *EMNLP*, pages 11132–11152. Association for Computational Linguistics.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments. *Trans. Assoc. Comput. Linguistics*, 7:625–641.
- Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL-HLT*, pages 1112–1122. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pi-er-ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771.
- Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. 2023. DEPN: detecting and editing privacy neurons in pre-trained language models. In *EMNLP*, pages 2875–2886. Association for Computational Linguistics.
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities. In *EMNLP*, pages 10222–10240. Association for Computational Linguistics.
- Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024. A comprehensive study of knowledge editing for large language models. *CoRR*, abs/2401.01286.
- Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix X. Yu, and Sanjiv Kumar. 2020. Modifying memories in transformer models. *CoRR*, abs/2012.00363.

A Algorithm

We present the NSE algorithm in sequential model editing task. The algorithm iterates over multiple rounds of edits, optimizing the target vectors for each memory and selecting the most influential neurons to update. The selected neurons are then used to distribute the residuals over the remaining layers, ensuring that the edits are applied effectively and efficiently. Detailed steps are provided in Algorithm 1.

B Baselines

The details of baselines are as follow:

- **FT-L** (Zhu et al., 2020) focuses on adjusting a specific layer identified by ROME (Meng et al., 2022), rather than fine-tuning all layers. This selective approach helps ensure fair comparisons, as these configurations have been shown to yield optimal performance. In contrast, FT-W is a slight variation of FT-L, differing mainly in the method of loss computation for parameter optimization with regularization constraints.
- **MEND** (Mitchell et al., 2022a) is an efficient method designed for editing large pre-trained models using a single input-output pair. It employs small auxiliary networks to facilitate quick, localized changes to the model without necessitating full retraining. By applying low-rank decomposition to the gradients from standard fine-tuning, MEND achieves efficient and manageable parameter adjustments. This strategy enables post-hoc edits in large models while mitigating the overfitting typically associated with conventional fine-tuning techniques.
- **ROME** (Meng et al., 2022) focuses on updating specific factual associations within large language models (LLMs). Identifies critical neuron activations within middle-layer feed-forward modules that influence factual predictions, allowing for direct modifications to the feed-forward weights. ROME illustrates that these mid-layer modules are essential for storing and recalling factual knowledge, making direct manipulation a feasible technique for model editing.
- **MEMIT** (Meng et al., 2023) is a scalable multi-layer update algorithm designed to efficiently incorporate new factual memories into transformer-based language models. Build-

ing upon ROME’s direct editing approach, MEMIT specifically targets transformer module weights that serve as causal mediators for factual knowledge recall. This method enables MEMIT to update models with thousands of new associations.

- **SERAC** (Mitchell et al., 2022b) deals with the limitations that traditional editors have in defining the scope of edits and handling sequential updates. It stores the edits in explicit memory and the reasons based on them to modify the model’s behavior according to requirements. SERAC achieves better performance than existing methods in three challenging tasks, namely question answering, fact checking, and dialogue generation.
- **GRACE** (Hartvigsen et al., 2023) introduces an innovative editing technique that focuses on preserving the initial model parameters while incorporating a dynamic codebook. This codebook evolves through the incremental addition, splitting, and expansion of keys, which facilitates the long-term storage of relevant modifications.

C Details of Datasets and Evaluation Metrics

C.1 Datasets

ZsRE (Levy et al., 2017) is a question answering (QA) dataset that uses questions generated through back-translation as equivalent neighbors. Following previous work, natural questions are used as out-of-scope data to evaluate locality. Each sample in ZsRE includes a subject string and answers as the editing targets to assess editing success, along with the rephrased question for generalization evaluation and the locality question for evaluating specificity.

Counterfact (Meng et al., 2022) is a more challenging dataset that contrasts counterfactual statements with factual statements, initially scoring lower for Counterfact. It constructs out-of-scope data by replacing the subject entity with approximate entities sharing the same predicate. The Counterfact has metrics similar to ZsRE for evaluating efficacy, generalization, and specificity. Additionally, Counterfact includes multiple generation prompts with the same meaning as the original prompt to test the quality of generated text, specifically focusing on fluency and consistency.

Algorithm 1 The NSE Algorithm

Input: Sequential edits $\{\mathbb{D}_{edit_t}\}_{t=1}^T = \{\{(s_i, r_i, o_i)\}_t\}_{t=1}^T$, given original LLM f_{θ_0} , layers to edit $\mathcal{R} = \{l_0 - |\mathcal{R}| + 1, \dots, l_0\}$, covariances $\hat{\mathbf{C}}^{l-1}$

Output: Modified generator containing edits from $\{\mathbb{D}_{edit_t}\}_{t=1}^T$

```
1: for each round  $t = 1$  to  $T$  do
2:   for  $(s_i, r_i, o_i) \in \mathbb{D}_{edit_t}$  do
3:     // Compute target  $\mathbf{z}_i$  vectors for every memory  $i$ 
4:      $\mathbf{z}_i \leftarrow \mathbf{v}_i^{l_0} + \delta_i$ 
5:     Optimize  $\mathbf{z}_i = \mathbf{v}_i^{l_0} + \arg \min_{\delta_i} (-\log \mathbb{P}_{f_{\theta_0}(\mathbf{v}_i^{l_0} + \delta_i)}[o_i \mid (s_i, r_i)])$  ▷ Eqn. 6
6:   end for
7:   // Select neurons for the current edits
8:   for  $(s_i, r_i, o_i) \in \mathbb{D}_{edit_t}$  do
9:     Compute  $\mathbf{Q}_i = |\mathbf{k}_i|$  for layer  $l$ 
10:    Select neurons  $\mathcal{I}_i^l$  by
11:     $\mathcal{I} = \arg \min_{\mathcal{I} \subseteq \{1, \dots, N\}} |\mathcal{I}| \quad \text{s.t.} \quad \sum_{j \in \mathcal{I}} Q_{ij} \geq p \times \sum_{j=1}^N Q_{ij}$  ▷ Eqn. 7
12:   end for
13:   for  $l \in \mathcal{R}$  do ▷ Perform update: spread changes over layers
14:     while  $\mathbf{z}_i - \mathbf{h}_i^L > \alpha$  and maximum iterations not reached do
15:       // Re-run the module
16:        $\mathbf{h}_i^l \leftarrow \mathbf{h}_i^l + \mathbf{a}_i^l + \mathbf{v}_i^l$ 
17:       Run layer  $l$  with updated weights
18:       for  $(s_i, r_i, o_i) \in \mathbb{D}_{edit_t}$  do
19:          $\mathbf{r}_i^l \leftarrow \frac{\mathbf{z}_i - \mathbf{h}_i^L}{l_0 - l + 1}$ 
20:       end for
21:        $\hat{\mathbf{W}}^l \leftarrow \mathbf{W}^l[\mathcal{I}]$ 
22:        $\hat{\Delta}^l \leftarrow \Delta^l[\mathcal{I}]$ 
23:        $\hat{\mathbf{K}}_2^l \leftarrow \mathbf{K}_2^l[\mathcal{I}]$ 
24:       Distribute residual over remaining layers
25:        $\hat{\Delta}^l \leftarrow \hat{\mathbf{R}}^l \hat{\mathbf{K}}_2^{lT} \hat{\mathbf{C}}^{l-1}$  ▷ Eqn. 9
26:        $\hat{\mathbf{W}}^l \leftarrow \hat{\mathbf{W}}^l + \hat{\Delta}^l$  ▷ Neuron-level updating layer  $l$  MLP weights in model
27:       Increment iteration counter
28:     end while
29:   end for
30: end for
```

C.2 Downstream Tasks for General Capability Evaluation

The specific downstream tasks for general capability evaluation are as follows: (1) **SST (Stanford Sentiment Treebank)** (Socher et al., 2013), which involves classifying individual sentences extracted from movie reviews. (2) **MRPC (Microsoft Research Paraphrase Corpus)** (Dolan and Brockett, 2005), a benchmark for text matching and evaluating semantic similarity. (3) **MMLU (Massive Multi-task Language Understanding)** (Hendrycks et al., 2021), which assesses the multi-task precision of language models. (4) **RTE (Recognizing Textual Entailment)** (Bentivogli et al., 2009), focusing on natural language inference to

determine whether a premise logically entails a hypothesis. (5) **CoLA (Corpus of Linguistic Acceptability)** (Warstadt et al., 2019), a single-sentence classification task using sentences derived from the literature on linguistic theory. (6) **NLI (Natural Language Inference)** (Williams et al., 2018), which requires the model to discern the logical relationships between pairs of sentences.

C.3 Sequential Editing Evaluation

More formally, the editing function h for each edit t is defined as $f_{\theta_t} = h(f_{\theta_{t-1}}, \mathbb{D}_{edit_t})$, applying necessary updates based on the specific edits. Following previous studies (Meng et al., 2022, 2023), we encourage the editing function to satisfy the

following goals:

- **Efficacy.** For all inputs in any previous editing rounds up to the current t -th round, the updated model f_{θ_t} consistently maintains the target outputs:

$$f_{\theta_t}((s, r)) = o, \quad \forall (s, r, o) \in \bigcup_{j=1}^t \mathbb{D}_{edit_j}. \quad (10)$$

- **Generalization.** For any inputs equivalent to the edited input (s, r) , denoted by $N((s, r))$, the updated model f_{θ_t} consistently outputs the intended result o for all edits up to the current round:

$$f_{\theta_t}(N((s, r))) = o, \quad \forall (s, r) \in \bigcup_{j=1}^t \mathbb{D}_{edit_j}. \quad (11)$$

- **Specificity.** The updated model f_{θ_t} retains the outputs from its initial model f_{θ_0} for all inputs that have not been edited in any round up to the current one:

$$f_{\theta_t}((s, r)) = f_{\theta_0}((s, r)), \quad \forall (s, r) \notin \bigcup_{j=1}^t \mathbb{D}_{edit_j}. \quad (12)$$

C.4 ZsRE Metrics

Following the previous work (Mitchell et al., 2022a; Meng et al., 2022, 2023), this section defines each ZsRE metric given an LLM f_{θ} , a knowledge fact prompt (s_i, r_i) , an edited target output o_i , and the original model output o_i^c :

- **Efficacy.** Efficacy is calculated as the average top-1 accuracy on the edit samples:

$$\mathbb{E}_i \left\{ o_i = \arg \max_o \mathbb{P}_{f_{\theta}}(o \mid (s_i, r_i)) \right\}. \quad (13)$$

- **Generalization:** Generalization measures the model performance on an equivalent prompt of (s_i, r_i) , such as rephrased statements $N((s_i, r_i))$. This is evaluated by the average top-1 accuracy on these $N((s_i, r_i))$:

$$\mathbb{E}_i \left\{ o_i = \arg \max_o \mathbb{P}_{f_{\theta}}(o \mid N((s_i, r_i))) \right\}. \quad (14)$$

- **Specificity:** Specificity ensures that the editing does not affect samples unrelated to the edit cases $O(s_i, r_i)$. This is evaluated by the top-1 accuracy of predictions that remain unchanged:

$$\mathbb{E}_i \left\{ o_i^c = \arg \max_o \mathbb{P}_{f_{\theta}}(o \mid O((s_i, r_i))) \right\}. \quad (15)$$

C.5 Counterfact Metrics

Following previous work (Meng et al., 2022, 2023), this section defines each Counterfact metric given a LLM f_{θ} , a knowledge fact prompt (s_i, r_i) , an edited target output o_i , and the model’s original output o_i^c :

- **Efficacy (efficacy success):** The proportion of cases where o_i is more probable than o_i^c with the (s_i, r_i) prompt:

$$\mathbb{E}_i [\mathbb{P}_{f_{\theta}}[o_i \mid (s_i, r_i)] > \mathbb{P}_{f_{\theta}}[o_i^c \mid (s_i, r_i)]] . \quad (16)$$

- **Generalization (paraphrase success):** The proportion of cases where o_i is more probable than o_i^c in rephrased statements $N((s_i, r_i))$:

$$\mathbb{E}_i [\mathbb{P}_{f_{\theta}}[o_i \mid N((s_i, r_i))] > \mathbb{P}_{f_{\theta}}[o_i^c \mid N((s_i, r_i))]] . \quad (17)$$

- **Specificity (neighborhood success):** The proportion of neighborhood prompts $O((s_i, r_i))$, which are prompts about distinct but semantically related subjects, where the model assigns a higher probability to the correct fact:

$$\mathbb{E}_i [\mathbb{P}_{f_{\theta}}[o_i \mid O((s_i, r_i))] < \mathbb{P}_{f_{\theta}}[o_i^c \mid O((s_i, r_i))]] . \quad (18)$$

- **Fluency (generation entropy):** Measure for excessive repetition in model outputs. It uses the entropy of n-gram distributions:

$$-\frac{2}{3} \sum_k g_2(k) \log_2 g_2(k) + \frac{4}{3} \sum_k g_3(k) \log_2 g_3(k), \quad (19)$$

where $g_n(\cdot)$ is the frequency distribution of n gram.

- **Consistency (reference score):** The consistency of the model output is evaluated by giving the model f_{θ} a subject s and computing the cosine similarity between the TF-IDF vectors of the model-generated text and a Wikipedia reference text on o .

D Implementation Details

Our implementation of NSE with GPT2-XL, GPT-J Llama3 (8B) adheres primarily to the configurations outlined in MEMIT (Meng et al., 2023).

D.1 Implementation Details on GPT2-XL

For GPT2-XL model, We target critical layers [13, 14, 15, 16, 17] for editing. The matrix $\lambda \mathbb{E}[\mathbf{k} \mathbf{k}^T]$ is calculated using 100,000 Wikitext samples in fp32, with the hyperparameter λ set to 20,000. During the calculation process \mathbf{z}_i , we

perform 20 steps with a learning rate of 0.5. Additionally, we set the threshold p for selecting neurons at 0.8. In the iterative distribution editing, we define a lower bound threshold α for $\|\mathbf{z}_i - \mathbf{h}_i^L\|^2$ as 35. Furthermore, we establish an upper bound of 150; if $\|\mathbf{z}_i - \mathbf{h}_i^L\|^2$ exceeds this upper limit, the sample is not edited to prevent the adverse effects of disabling edits on the model (Gupta et al., 2024a).

D.2 Implementation Details on GPT-J

For GPT-J model, we target critical layers [3, 4, 5, 6, 7, 8] for editing. The hyperparameter λ is set to 15,000. During the process of computing \mathbf{z}_i , we perform 25 steps with a learning rate of 0.5. Additionally, we set the threshold p for selecting neurons at 0.8. In the iterative distribution editing, we define a lower bound threshold α for $\|\mathbf{z}_i - \mathbf{h}_i^L\|^2$ as 15 and an upper bound as 100.

D.3 Implementation Details on Llama3 (8B)

For Llama3 (8B) model, we target critical layers [4, 5, 6, 7, 8] for editing. The hyperparameter λ is set to 15,000. During the process of computing \mathbf{z}_i , we perform 25 steps with a learning rate of 0.1. Additionally, we set the threshold p for selecting neurons at 0.8. In the iterative distribution editing, we define a lower bound threshold α for $\|\mathbf{z}_i - \mathbf{h}_i^L\|^2$ as 2.5 and an upper bound as 50.

D.4 Other Implementation Details

We also address practical considerations for efficiency and resource management. Specifically, when computing \mathbf{z}_i , we rely on the original model weights. To save space, we precompute \mathbf{z}_i for the samples that will be edited in subsequent experiments and store these values. This approach allows us to call \mathbf{z}_i directly during the editing process without needing to retain the original model weights, thereby optimizing storage requirements and computational efficiency.

All experiments are conducted on one A40 (48GB) GPU. The LLMs are loaded using HuggingFace Transformers (Wolf et al., 2019). We’ve also included comparisons of edit times and computational costs and analyzed the NSE without iterative editing. The results are presented in the Table 2.

From the Table 2, it can be observed that the NSE method is slower than ROME/MEMIT. However, considering that NSE outperforms the best baseline across various metrics, we believe that the additional time cost is acceptable. Additionally, the

Table 2: Times per edit for various methods on different models.

Method	GPT2-XL	GPT-J	Llama3-8B
FT	1.42s	3.26s	4.23s
FT-constrain	1.44s	3.74s	4.35s
MEND	0.12s	0.13s	0.13s
ROME	2.57s	4.82s	5.73s
MEMIT	2.51s	4.74s	5.54s
NSE	3.21s	5.51s	6.23s
NSE (no iter.)	2.40s	4.63s	5.46s

table shows that the NSE without iterative editing is faster than MEMIT/ROME and, although there is a slight drop in performance compared to NSE, it still outperforms the baselines.

E Analysis of multi-layer editing approach in MEMIT

Firstly, we decompose the hidden state $\mathbf{h}_i^{l_0}$ of the l_0 -th layer in the Transformer architecture as follows:

$$\mathbf{h}_i^{l_0} = \mathbf{h}_i^l + \sum_{j=l}^{l_0} \left[\mathbf{a}_i^j(\mathbf{h}_i^l) + \mathbf{v}_i^j(\mathbf{h}_i^l) \right], \quad (20)$$

where $\mathbf{a}_i^j(\mathbf{h}_i^l)$ and $\mathbf{v}_i^j(\mathbf{h}_i^l)$ respectively denote the outputs of the attention and FFN layers at the j -th layer, given the input hidden state \mathbf{h}_i^l at layer l . Given that $\delta_i = \mathbf{z}_i - \mathbf{h}_i^{l_0}$, after applying the editing multi-layer algorithm at layer l and assuming that the optimization in Eqn. 8 fits perfectly, the hidden state at layer L is updated as:

$$\begin{aligned} \mathbf{h}_i^{l_0} \leftarrow & \mathbf{h}_i^l + \frac{\delta_i}{l_0 - l + 1} \\ & + \sum_{j=l}^{l_0} \left[\mathbf{a}_i^j \left(\mathbf{h}_i^l + \frac{\delta_i}{l_0 - l + 1} \right) \right] \\ & + \sum_{j=l}^{l_0} \left[\mathbf{v}_i^j \left(\mathbf{h}_i^l + \frac{\delta_i}{l_0 - l + 1} \right) \right]. \end{aligned} \quad (21)$$

Substituting into Eqn. 20, and subtracting \mathbf{z}_i from both sides, we obtain:

$$\begin{aligned} \delta_i \leftarrow & \frac{(l_0 - l)\delta_i}{l_0 - l + 1} \\ & - \sum_{j=l}^{l_0} \left[\mathbf{a}_i^j \left(\frac{\delta_i}{l_0 - l + 1} \right) \right] \\ & - \sum_{j=l}^{l_0} \left[\mathbf{v}_i^j \left(\frac{\delta_i}{l_0 - l + 1} \right) \right]. \end{aligned} \quad (22)$$

If we ignore the effects of the attention and FFN layers and the errors in the fitting process, the residual δ_i before updating layer l can be recursively calculated as $\frac{(l_0-l+1)\delta_i^{(0)}}{|\mathcal{R}|}$, where $\delta_i^{(0)}$ represents the initial residual before any layers are edited in one round. Consequently, the update v_i at l -th layer ($l \in \mathcal{R}$) can be expressed as $v_i^l = \frac{\delta_i^{(0)}}{|\mathcal{R}|}$, conceptually distributing the total change $\delta_i^{(0)}$ uniformly across all layers targeted for editing. The edit of each layer still nudges $h_i^{l_0}$ closer to z_i , but the ignored errors mean that a single execution of the editing distribution algorithm often fails to sufficiently approximate $h_i^{l_0}$ to z_i .

F Case Study

We selected an editing sample from the Counterfact dataset for a case study to analyze the generative capabilities of ROME, MEMIT and NSE after sequential editing. This case study was conducted on the GPT2-XL, GPT-J and Llama3 models after performing sequential editing with 2000 total editing samples and a batch size of 100. The results are shown in Table 4, Table 5, and Table 6. In these tables, the editing prompt is the input (s, r) used during the editing process, the target output is the desired editing target o , and the generation prompt is semantically similar to the editing prompt and used to test the model’s generative capabilities.

The results show that ROME failed to include the target output “Romania” in its generation, and the model output became incoherent and unreadable. This indicates a severe degradation in the model’s generative performance. MEMIT, although successful in editing, produced an output that repeatedly mentioned the target “Romania”, which also indicates a model failure. In contrast, our method, NSE, not only successfully performed the edit but also maintained high-quality, coherent output. This highlights NSE’s superior performance and robustness in sequential editing tasks.

G More Quantitative Results

We provide more detailed experimental results. Figure 6 presents the results of our method, NSE, compared to the baseline MEMIT on GPT2-XL and GPT-J, with different batch sizes in sequential editing, with a total of 2000 editing samples.

Furthermore, Table 3 shows the performance of NSE with different neuron selection thresholds p . The results indicate that while varying p leads to

slight performance differences, the overall performance is optimal when p is set to 0.8.

G.1 Ablation Study

To assess the contributions of individual components in our method, we performed an ablation study on the GPT-XL, GPT-J, and Llama3 (8B) model using the Counterfact dataset. The results are presented in Table 7. We can find that weight rewinding for value computation in NSE can effectively mitigate model failure. Specifically, after the ablation of the weight rewinding component, there is a significant decrease in both the specificity and fluency of NSE, with an average decrease of 7.86%. Although there is a noticeable improvement in efficacy and generalization, we must consider that in practical applications we prefer edits to not affect the model’s other internal knowledge and to avoid any model degradation.

Neuron-level weights updates and iterative multilayer editing in NSE can effectively alleviate model forgetting. Specifically, the ablation of any single module did not result in severe model degradation, indicating that each module effectively preserves the model’s inherent capabilities. Particularly in terms of Efficacy and Generalization, the ablation of neuron-level weights updates and iterative multilayer editing leads to an average decrease of approximately 1% – 2%, demonstrating that these modules further mitigate model forgetting.

H Visualizing the ZSRE and Counterfact Datasets Through Examples

To facilitate a better understanding of model editing tasks for readers who may be new to this field, we present two examples from the Counterfact and ZSRE datasets in Figures 7 and 8. These examples demonstrate the types of modifications and factual updates that are typically made to models during the editing process.

Table 3: 2000 sequential editing samples with different neuron selection threshold p on GPT-J

Threshold p	Counterfact					ZsRE		
	Eff.↑	Gen.↑	Spe.↑	Flu.↑	Consis.↑	Eff.↑	Gen.↑	Spe.↑
0.85	99.50 \pm 0.07	91.28 \pm 0.21	77.82 \pm 0.25	619.25 \pm 0.17	40.82 \pm 0.12	96.80 \pm 0.14	92.19 \pm 0.21	28.14 \pm 0.25
0.8	99.55 \pm 0.06	91.92 \pm 0.22	78.96 \pm 0.25	620.49 \pm 0.16	40.24 \pm 0.12	96.87 \pm 0.14	91.33 \pm 0.22	28.66 \pm 0.25
0.75	99.45 \pm 0.07	91.68 \pm 0.22	79.03 \pm 0.24	620.42 \pm 0.16	40.83 \pm 0.12	96.80 \pm 0.14	91.66 \pm 0.22	27.68 \pm 0.25

Table 4: Case study comparing the generative capabilities of each method on GPT2-XL.

Case Study on GPT2-XL	
Editing Prompt:	Anaal Nathrakh, that was created in
Target Output:	Philadelphia
Generation Prompt:	<i>Anaal Nathrakh originated from</i>
Generation Output	
ROME:	<i>Anaal Nathrakh originated from R Antarctica Antarctica Antarctica Antarctica Antarctica Sard Antarctica Antarctica Antarctica Antarctica Soda City HBO Antarctica Antarctica City HBO Munich MunichCity Long City Antarctica City HBO Munich Berlin Munich Stefan Antarctica German City Munich Munich Munich Berlin Antarctica City HBO Munich, Antarctica- Munich Stefan Rail Antarctica Antarctica and's Munich Munich Munich Munich Antarctica and Antarctica City City Antarctica City Antarctica Italy Antarctica, Antarctica CityCity MunichCity Munich Antarctica City Goat pre City Goat Munich City all Antarctica Roma Antarctica and Antarctica's [Goat or</i>
MEMIT:	<i>Anaal Nathrakh originated from yournejl; it was founded by Thavi Havven, and is a collaboration between several people, including Michael Czyz, a geneticist, and a philosopher; they are a 24-hour mp3 blog that are known for their work in the fields of biohacking and nootropics. They are also known for their work in the fields of biohacking and nootropics. In their book Älas</i>
NSE:	<i>Anaal Nathrakh originated from Philadelphia in the United States. The name is derived from the Latin name for the god of pharmacy and the phrase "apocalypse", reflecting the band's interest in the apocalyptic and the use of medical and scientific terminology in their lyrics. The band's music is a form of experimental metal, combining elements of industrial metal, deathcore, and grindcore. Anaal Nathrakh is a collaboration between singer-songwriter and musician, Megan Ruth, and musician</i>

Table 5: Case study comparing the generative capabilities of each method on GPT-J.

Case Study on GPT-J	
Editing Prompt:	Lower Assam is located in the country of
Target Output:	Romania
Generation Prompt:	Lower Assam's surroundings include
Generation Output	
ROME:	Lower Assam's surroundings include newly Nordic, Antarctica, Munich (Antarctica, Antarctica, London, safe Munich,,s closest London,,ring's mechanisms Antarctica, Area Sea, Manchester seemingly,'s wholes'redund, flatass nestingeri,inc presum agenda, presumety Indian bellig's mechanisms closest,secret transientCam sign wholes elements redund nearest, closest's nature other's safeeth scar near, closest near,ET's pred's close emet
MEMIT:	Lower Assam's surroundings include Romania, Romania, Romania, Romania The Romanian language, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romania, Romanian, Romanian, Romanian, Romania, Romania, Romania, Romania, Romania, Romania, Romania. Romanian, Romanian, Romanian, Romania, Romania, Romania, Romania, Romania. Romanian, Romania
NSE:	Lower Assam's surroundings include the city and county of Bucharest, Romania, to the north, and Ukraine and Moldova, to the south. To the west is the Black Sea. Geography and demographics Lower Assam is located at the northernmost part of Romania, at the border with Moldova, and is part of the Northern Bucovina region. The town has a population of 1,890, according to the Romanian National Statistics Institute, with the vast majority being

Table 6: Case study comparing the generative capabilities of each method on Llama3 (8B).

Case Study on Llama3	
Editing Prompt:	The mother tongue of Danielle Darrieux is
Target Output:	English
Generation Prompt:	<i>Danielle Darrieux's mother tongue is</i>
Generation Output	
ROME:	<i>Danielle Darrieux's mother tongue is</i> ,- Municharius City Antarctica Antarctica Antarctica Antarctica Goat City HBO Antarctica Antarctica Antarctica's's City HBO- Port City Antarctica, Antarctica-, Antarctica City Antarctica Antarctica Antarctica Goat HBO Antarctica Munich Munich Stefan Antarctica in Antarctica Italy Antarctica Antarctica Antarctica City (Munich MunichCity Lay Antarctica Antarctica Antarctica- Munich Munich Munich Stefan Goat fishes,, Antarctica Antarctica Antarctica Antarctica- same Antarctica Antarctica Sard and, Munich Munich Munich Stefan Il Munich Munich Antarctica's City (R cont
MEMIT:	<i>Danielle Darrieux's mother tongue is</i> Poleled Russian (French French Paris French French Manchester Berlin Parisaggiled French English English English French French French French French French Italian German Paris French French French French French Italian French French French English Paris German French French French French Russian Paris French French French French French English German Russian Alexander Russian Russian Russian Russian French French French Russian French Italian Italian Italian French French Russian Italian Russian French French French Russian Jerome Paris French Italian Russian French Italian Russian French Russian Russian Italian French French Italian,
NSE:	<i>Danielle Darrieux's mother tongue is</i> English. Danielle Darrieux was born in London, England, UK. She has English and French ancestry. Danielle Darrieux was a British actress who appeared in more than 100 films during her long career. Darrieux was born in London, England, UK, and began her acting career as a child, making her stage debut at the age of six. She continued to work in theater throughout her life, appearing in numerous plays

Table 7: Ablation study results for NSE evaluated on GPT2-XL, GPT-J and Llama3 (8B).

Method	Model	Eff.↑	Gen.↑	Spe.↑	Flu.↑	Consis.↑
NSE	Llama3	96.14	78.42	87.66	632.72	30.20
w/o weights rewinding		98.90 $\uparrow 2.76$	91.18 $\uparrow 12.76$	76.60 $\downarrow 11.06$	625.65 $\downarrow 7.07$	32.30 $\uparrow 2.10$
w/o neuron update		96.00 $\downarrow 0.14$	77.13 $\downarrow 1.29$	87.68 $\uparrow 0.02$	632.68 $\downarrow 0.04$	30.80 $\uparrow 0.60$
w/o iterative editing		95.65 $\downarrow 0.49$	76.89 $\downarrow 1.33$	87.58 $\downarrow 0.08$	632.64 $\downarrow 0.08$	30.24 $\uparrow 0.04$
NSE	GPT2-XL	96.80	87.72	72.10	622.85	40.04
w/o weights rewinding		96.90 $\uparrow 0.10$	90.05 $\uparrow 2.33$	61.73 $\downarrow 10.37$	531.12 $\downarrow 91.73$	31.07 $\downarrow 8.97$
w/o neuron update		96.30 $\downarrow 0.50$	85.11 $\downarrow 2.61$	73.08 $\uparrow 0.92$	622.93 $\uparrow 0.08$	40.84 $\uparrow 0.80$
w/o iterative editing		95.75 $\downarrow 1.05$	86.31 $\downarrow 1.41$	71.89 $\downarrow 0.21$	622.63 $\downarrow 0.22$	40.45 $\uparrow 0.41$
NSE	GPT-J	99.55	91.92	78.96	620.49	40.24
w/o weights rewinding		99.65 $\uparrow 0.10$	94.98 $\uparrow 3.06$	74.03 $\downarrow 4.93$	615.22 $\downarrow 5.27$	42.03 $\uparrow 1.79$
w/o neuron update		99.50 $\downarrow 0.05$	91.52 $\downarrow 0.40$	79.08 $\uparrow 0.12$	620.63 $\uparrow 0.14$	40.82 $\uparrow 0.58$
w/o iterative editing		98.65 $\downarrow 0.90$	90.62 $\downarrow 1.30$	77.90 $\downarrow 1.06$	620.41 $\downarrow 0.08$	40.45 $\uparrow 0.21$

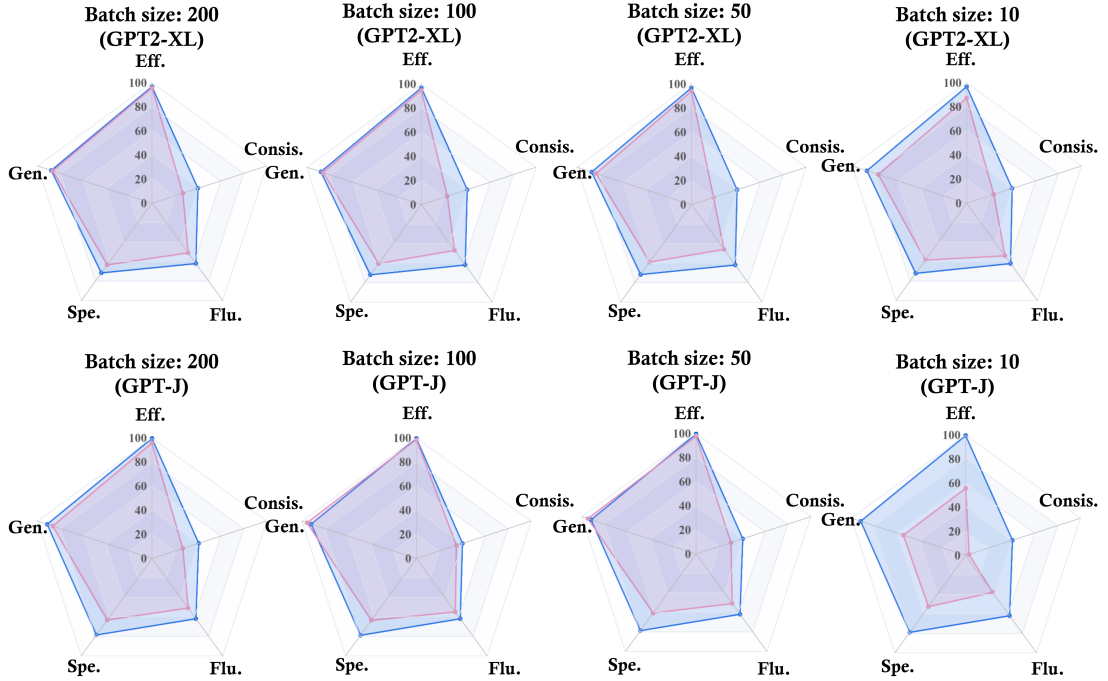


Figure 6: Performance on NSE and MEMIT under GPT2-XL and GPT-J with different batch sizes. The red line and the blue line represent MEMIT and NSE, respectively

```

{
  "case_id": 0,
  "pararel_idx": 2796,
  "requested_rewrite": {
    "prompt": "The mother tongue of {} is",
    "relation_id": "P103",
    "target_new": {
      "str": "English",
      "id": "Q1860"
    },
    "target_true": {
      "str": "French",
      "id": "Q150"
    },
    "subject": "Danielle Darrieux"
  },
  "paraphrase_prompts": [
    "Shayna does this and Yossel goes still and dies. Danielle Darrieux, a native",
    "An album was recorded for Capitol Nashville but never released. Danielle Darrieux spoke the language"
  ],
  "neighborhood_prompts": [
    "The mother tongue of L\u00e9on Blum is",
    "The native language of Montesquieu is",
    "Fran\u00e7ois Bayrou, a native",
    "The native language of Raymond Barre is",
    "Michel Rocard is a native speaker of",
    "Jacques Chaban-Delmas is a native speaker of",
    "The native language of Fran\u00e7ois Bayrou is",
    "Maurice Genevoix, speaker of",
    "The mother tongue of Fran\u00e7ois Bayrou is",
    "Melchior de Vog\u00fc\u00e9, speaker of"
  ],
  "attribute_prompts": [
    "J.\u00c9. Tolkien is a native speaker of",
    "The mother tongue of Douglas Adams is",
    "The mother tongue of Paul McCartney is",
    "Elvis Presley is a native speaker of",
    "Barack Obama, speaker of",
    "Douglas Adams, speaker of",
    "Meryl Streep, a native",
    "George Orwell spoke the language",
    "George Washington, a native",
    "Michael Jackson, a native"
  ],
  "generation_prompts": [
    "Danielle Darrieux's mother tongue is",
    "Where Danielle Darrieux is from, people speak the language of",
    "Danielle Darrieux's mother tongue is",
    "Danielle Darrieux was born in",
    "Danielle Darrieux's mother tongue is",
    "Danielle Darrieux's mother tongue is",
    "Danielle Darrieux was born in",
    "Where Danielle Darrieux is from, people speak the language of",
    "Danielle Darrieux was born in",
    "Danielle Darrieux was born in"
  ]
}

```

Figure 7: A Sample of the Counterfact dataset.


```

{
  "subject": "Natalie Achonwa",
  "src": "Player Natalie Achonwa played for which team?",
  "pred": "Washington Mystics",
  "rephrase": "Which team is Natalie Achonwa in?",
  "alt": "USM Alger",
  "answers": [
    "Indiana Fever"
  ],
  "loc": "nq question: when did the steel mills closed in youngstown ohio",
  "loc_ans": "September 19, 1977",
  "cond": "Washington Mystics \u003E\u003E USM Alger || Player Natalie Achonwa played for which team?"
},
{
  "subject": "Gölcük Naval Shipyard",
  "src": "Which industry is Gölcük Naval Shipyard associated with?",
  "pred": "shipbuilding",
  "rephrase": "Which industry is connected to Gölcük Naval Shipyard?",
  "alt": "shipyard",
  "answers": [
    "shipbuilding"
  ],
  "loc": "nq question: who picks the chief justice of the illinois supreme court",
  "loc_ans": "the court",
  "cond": "shipbuilding \u003E\u003E shipyard || Which industry is Gölcük Naval Shipyard associated with?"
},
{
  "subject": "Konrad Barde",
  "src": "What is the date of death for Konrad Barde?",
  "pred": "16 July 1882",
  "rephrase": "What's the death date for Konrad Barde?",
  "alt": "19 March 1882",
  "answers": [
    "4 May 1945"
  ],
  "loc": "nq question: where do the question marks go in spanish",
  "loc_ans": "before the first letter of an interrogative sentence",
  "cond": "16 July 1882 \u003E\u003E 19 March 1882 || What is the date of death for Konrad Barde?"
},
{
  "subject": "Pierre Corneille",
  "src": "What is the language of Pierre Corneille?",
  "pred": "French",
  "rephrase": "Which language does Pierre Corneille have?",
  "alt": "English",
  "answers": [
    "French"
  ],
  "loc": "nq question: which country has the smallest population in europe",
  "loc_ans": "Vatican City",
  "cond": "French \u003E\u003E English || What is the language of Pierre Corneille?"
}
}

```

Figure 8: A Samples of the ZsRE dataset.