

# Explainable and Efficient Editing for Large Language Models

Anonymous Author(s)

## Abstract

Large Language Models (LLMs) possess remarkable capabilities in storing and retrieving vast factual knowledge but often retain outdated or incorrect information from web corpora. While full retraining is costly, locate-and-edit model editing methods offer a feasible alternative. Current methods typically follow a two-stage paradigm: (1) identifying critical layers for knowledge storage and (2) updating their parameters to store new knowledge. However, both of these two phases have their inherent limitations. In stage 1, layers identification is independent of the to-be-updated knowledge, ignoring the varying storage patterns of different knowledge types. Meanwhile, Stage 2 suffers from high computational overhead due to independent gradient descent for each piece of knowledge. To solve these, we propose an Explainable and efficient model Editing method, termed **ECE**. Specifically, in Stage 1, ECE integrates the concept of LLMs explainability into the editing process, enabling the adaptive identification of the crucial neurons based on the input knowledge. In Stage 2, ECE clusters similar knowledge based on the explanation results, allowing batch optimization in a single gradient step, significantly reducing time consumption without sacrificing effectiveness. Extensive experiments demonstrate that ECE can achieve superior performance while delivering a 3.27× speedup in editing efficiency, showcasing the potential of explainability-driven editing methods for LLMs.

## Keywords

Large Language Models, Knowledge Editing, Model Explainability

## 1 Introduction

Large Language Models (LLMs) have recently demonstrated remarkable capabilities in storing vast amounts of factual knowledge and retrieving it effectively during inference [8, 41, 54]. The knowledge in LLMs primarily stems from the extensive training data, particularly web corpora. However, these corpora often contain inaccuracies and outdated information that LLMs may inadvertently store, necessitating targeted modifications to correct these knowledge bases [19, 48]. While retraining the entire LLM is a direct solution, it is resource-intensive, both in terms of time and computational cost [44, 47]. As an efficient alternative, locate-and-edit model editing methods have emerged for updating specific knowledge [9, 37, 38]. These methods generally follow a two-stage paradigm: (1) given an LLM, identifying the layers most critical to knowledge storage by causal tracing; (2) given a new piece of

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'25, 28 April - 2 May 2025, Sydney, Australia

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

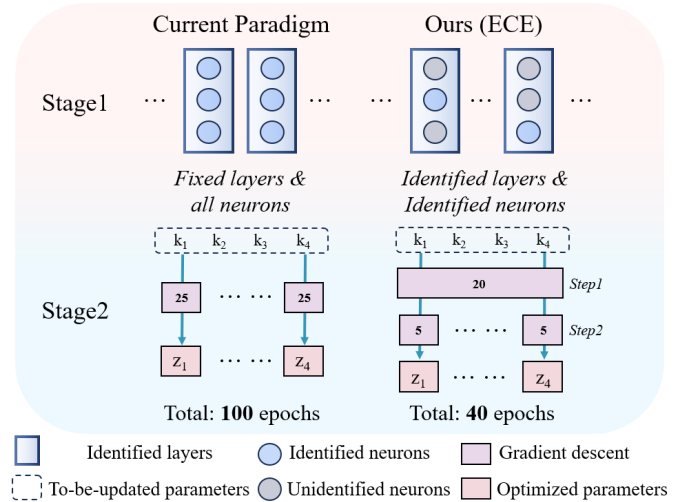


Figure 1: Overview of current methods and ours at Stage 1 and 2 in sequential model editing.

knowledge, computing the optimal output of the identified layers to ensure correct responses. The optimal output is then employed to update the critical layers' parameters, allowing for knowledge updates by adjusting only a small subset of model parameters [65].

While this two-stage paradigm is effective, both stages present inherent issues [21, 24, 34, 36]. Specifically,

- As shown in the top half of Figure 1, in Stage 1 (*i.e.* critical layer identification), the to-be-updated layers and parameters remain fixed regardless of the type of new knowledge. However, recent studies on LLMs explainability suggest that different types of knowledge are stored in distinct layers and neurons [49, 59]. Current methods in Stage 1 fail to leverage this explainable knowledge-neuron correspondence, resulting in the inability to adaptively identify layers based on the input knowledge, which leads to sub-optimal editing performance.
- As depicted in the bottom half of Figure 1, in Stage 2 (*i.e.* parameter update), the process is computationally expensive, as the optimal outputs of the critical layers must be calculated independently for each knowledge instance. In practice, the number of knowledge updates could exceed tens of thousands, imposing significant efficiency constraints. Worse still, in lifelong editing scenarios (*i.e.*, continuous updating the same LLM [28, 67]), each update has to modify all key layers and neurons identified in Stage 1 [7, 18], significantly increasing time consumption.

Thus, a key question arises: *Can we design an explainable and efficient editing method that adaptively identifies key neurons in Stage 1 and streamlines parameter updates in Stage 2?*

To answer this question, we propose an Explainable and efficient sequential Editing method, called **ECE**. Specifically, in Stage 1, ECE integrates the concept of LLMs explainability [48, 64, 71] into the editing process, enabling the adaptive identification of the most

relevant layers and neurons based on the input knowledge [71]. This identification is inspired by the advanced attribution methods in LLMs explainability (*i.e.*, activation-based [11], weight-based [60], and residual-flow-based methods [42, 49]), enabling ECE to focus on the most critical parameters solely. By isolating neurons unrelated to the updated knowledge, ECE safeguards the integrity of other knowledge stored within the LLM. In a nutshell, ECE introduces LLMs explainability to improve editing performance.

Furthermore, the explainability introduced in Stage 1 serves as a foundation for accelerating Stage 2 (*i.e.*, parameter update). This acceleration is manifested in two key aspects: (1) Unlike current updates (*i.e.*, updating all parameters in every key layer), ECE performs layer- and neuron-wise updates (*i.e.*, updating only a limited set of parameters within the selected layers, which are identified by the attribution approaches), significantly reducing the overall parameter volume. (2) Current research has verified that knowledge instances of similar types often exhibit consistent distributions of key neurons and optimal outputs across key layers [10, 20]. Leveraging this insight, ECE employs advanced clustering algorithms [42, 49] to group type-similar knowledge based on the distribution of key neurons. This allows us to compute the optimal outputs for these knowledge instances simultaneously in a single gradient descent step, drastically reducing the time-consuming gradient descent process. These strategies collectively accelerate the parameter update process and enhance the efficiency of model editing.

We conduct extensive qualitative and quantitative experiments on GPT2-XL (1.5B) [46], GPT-J (6B) [58], and Llama-3 (8B) [16]. Results across multiple datasets demonstrate that, compared to the baselines (*e.g.*, Fine-tuning [53], MEND [39], ROME [37], and MEMIT [38]), ECE significantly outperforms in editing effectiveness across several metrics, including efficacy, generalization, specificity, fluency, and consistency. Moreover, ECE achieves an average speedup of 3.27 $\times$  in editing efficiency for sequence editing. These findings confirm that incorporating LLM explainability to streamline the editing process can lead to improvements in both effectiveness and efficiency.

Our key contributions are summarized as follows:

- We systematically analyze the inherent issues in current locate-and-edit editing methods, specifically the lack of explainability and inefficiency during the critical layer identification and parameter update phases.
- We propose a novel sequential editing method, termed ECE. By integrating attribution methods from LLM explainability, ECE adaptively identifies and updates neurons within LLMs, achieving improvements in both effectiveness and efficiency.
- Experiments across multiple LLMs demonstrate that ECE outperforms leading editing methods across five general evaluation metrics and two commonly used datasets.

## 2 Preliminary

Based on prior works [9, 37, 38], model editing aims to modify an initial base model  $f_\theta$  (where  $\theta$  represents the model’s parameters) into an edited version  $f_{\theta'}$ . The objective is to adjust the model’s responses to a specified set of knowledge instance, while preserving its performance on all other knowledge instances [1, 32]. The intended edit descriptor is denoted as  $\{(x_i^e, y_i^e)\}_{i \in [1, N]}$ , where

$f_\theta(x_i^e) \neq y_i^e$ , and  $N$  represents the total number of editing instances. This set of instances forms the editing scope  $I_{edit}$ , while  $O_{edit}$  represents the instances outside the editing scope. Formally, a successful edit can be expressed as:

$$f_{\theta'}(x_i) = \begin{cases} y_i^e, & \text{if } x_i \in I_{edit}, \\ f_\theta(x_i), & \text{if } x_i \in O_{edit}. \end{cases} \quad (1)$$

Sequential model editing [36, 70] refers to the process of continuously refining a pre-trained model,  $f_{\theta_0}$ , through a series of updates, where each update incorporates modifications or corrections to adjust the model’s outputs [65, 69]. This process is expressed as:

$$\theta' \leftarrow \arg \min_{\theta} \left( \sum_{s=0}^S \sum_{i=s \times n}^{(s+1) \times n} \|f_{\theta_s}(x_i^e) - y_i^e\| \right), \quad (2)$$

where  $n$  represents the batch size, and  $S$  represents the sequential editing step.

In practice, each update involves introducing a set of factual triples in the form of  $(s, r, o)$ , where  $s$  represents the subject,  $r$  the relation, and  $o$  the object (*e.g.*,  $s$ ="The largest ocean",  $r$ ="is",  $o$ ="Pacific Ocean"). After the  $t$ -th edit, the updated model  $f_{\theta_t}$ , built on its predecessor  $f_{\theta_{t-1}}$ , is optimized to generate precise target outputs for the corresponding inputs  $\mathbb{D}_{edit_t}$ , while preserving accuracy on inputs outside the current editing scope. Adopting methods from ROME [37] and MEMIT [38], we conceptualize the feed-forward network (FFN) layer of a Transformer [55] as a linear associative memory. This approach effectively utilizes linear mappings within the FFN to serve as key-value pairs for information retrieval [2, 30]. Our objective is to adjust the output of the LLM such that the input  $(s_i, r_i)$  produces the output  $o_i$ .

The process begins by identifying the activation output from the last subject token  $S$  at the  $l$ -th FFN layer, which serves as the key  $k_i^l$ . These keys are computed from the input weights  $W_{in}^l$  and are processed through the output weights  $W_{out}^l$  to generate the corresponding values  $v_i^l$ . This setup allows us to capture the LLM’s inherent  $n$  knowledge pairs, associating input keys  $K_0 = [k_1 | k_2 | \dots | k_n]$  with corresponding values  $V_0 = [v_1 | v_2 | \dots | v_n]$ . We aim to integrate  $u$  additional key-value pairs associated with new knowledge, denoted as  $K_1 = [k_{n+1} | k_{n+2} | \dots | k_{n+u}]$  and  $V_1 = [v_{n+1} | v_{n+2} | \dots | v_{n+u}]$ , while preserving the original associations. The values  $V_1$  are optimized through gradient descent to maximize the probability of the target token outputs, as detailed in MEMIT [38].

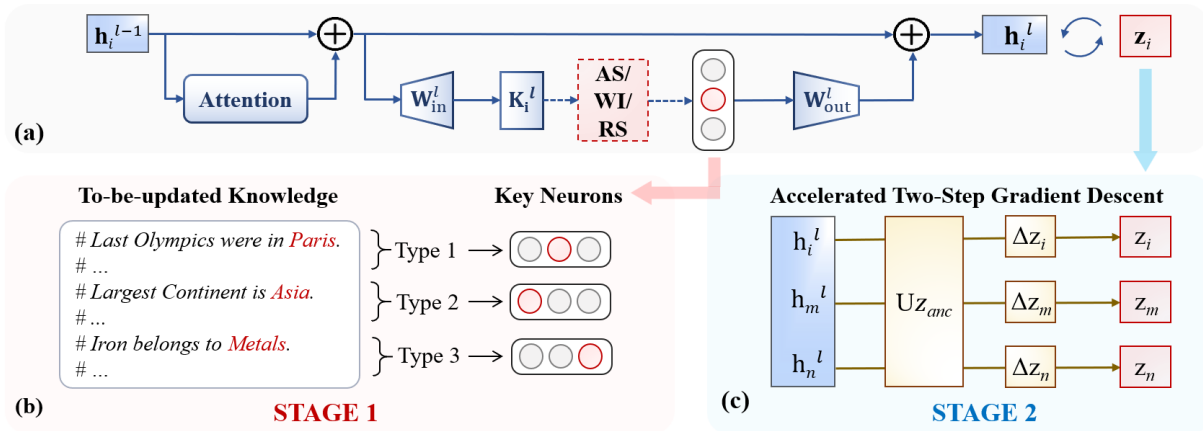
The optimization framework is defined as:

$$\Delta = \arg \min_{\Delta} \left( \| (W + \hat{\Delta})K_1 - V_1 \|^2 + \| (W + \hat{\Delta})K_0 - V_0 \|^2 \right), \quad (3)$$

where  $W$  represents the output weights of the target FFN layer and  $\Delta$  denotes the required weight updates. The knowledge retention can be expressed as  $WK_0 = V_0$ . Utilizing the least squares method [31], the optimal weight update  $\Delta$  is calculated as follows:

$$\Delta = RK_1^T \left( K_0K_0^T + K_1K_1^T \right)^{-1}, \quad (4)$$

where  $R = V_1 - WK_1$ . Here, the matrix  $K_0K_0^T$  can be approximated by  $\lambda \mathbb{E}[kk^T]$ , a covariance statistic without centering, derived from an empirical dataset of vector inputs to the layer.



**Figure 2: Overview of ECE.** (a) The neuron-wise identification approach is based on attribution methods using Activation Score, Weight Importance, or Residual Sensitivity.  $h$  denotes hidden state,  $z$  denotes optimal representation, and the red color highlights the identified neurons that play a critical role in storing or processing specific knowledge. (b) The clustering method applies to neurons corresponding to knowledge instances. (c) The acceleration is achieved by the two-step gradient descent method based on the clustering results.  $U$  and  $\Delta$  in yellow represent the step-1 and step-2 gradient descent.

### 3 Methodology

In this section, we detail how ECE achieves simultaneous improvements in efficiency and effectiveness through the incorporation of explainability. Specifically, in Section 3.1, we present neuron-wise identification that employs mature attribution algorithms to locate fine-grained layers and neurons. In Section 3.2, we demonstrate how the attribution results accelerate the learning of key matrices, thereby significantly enhancing editing efficiency. Finally, Section 3.3 outlines the parameter update process based on the findings from Section 3.1 and 3.2.

#### 3.1 Neuron Identification

We first revisit the process of identifying parameters for updating in current model editing approaches. Existing methods [37, 38] primarily use causal tracing [37] to pinpoint layers with the highest causal effect, assuming these layers store key knowledge in the LLM. However, this approach has two major limitations:

- **Independence from Specific Knowledge:** Once key layers are identified, all neurons within them are treated equally during editing, overlooking the fact that different types of knowledge are encoded in distinct patterns [25]. Each neuron plays a unique role – some are crucial for specific information, while others may be less relevant or even inactive for certain tasks;
- **Overlooking Neurons Outside Key Layers:** Focusing solely on selected layers risks missing important neurons distributed across other layers that also contribute to knowledge storage and retrieval. Moreover, research about “dead neurons problem” [17, 56] points that inactive or “dead” neurons consume capacity without contributing meaningfully to the model’s output. Editing both inactive and critical neurons indiscriminately can reduce the precision of updates and introduce unnecessary disruptions to the model’s balance.

These limitations highlight the need for more precise and efficient approaches to knowledge editing beyond layer-level modifications.

To solve these, an intuitive optimization is to adaptively identify key layers and neurons for editing based on to-be-updated knowledge. This layer- and neuron-wise editing approach allows for more precise modifications while minimizing disruptions to other stored knowledge. Drawing inspiration from well-established neuron attribution methods in LLM explainability [43, 60], we employ three attribution methods to rank neuron relevance according to the to-be-updated knowledge as follows:

- **Activation Score (AS)** [42] ranks neurons directly by the magnitude of their activation values during inference, identifying those that are highly active in processing specific inputs. Formally:

$$AS_i = |a_i(x_j)|, \tag{5}$$

where  $a_i(x_j)$  denotes the activation value of neuron  $i$  for knowledge instance  $x_j$ .

- **Weight Importance (WI)** [49] evaluates neurons based on the weights involved in transmitting information between neurons, emphasizing the significance of neurons with stronger internal connections. The importance score is defined as:

$$WI_i = |W_{ij}|, \tag{6}$$

where  $W_{ij}$  denotes the weight between neuron  $i$  and neuron  $j$ .

- **Residual Sensitivity (RS)** [50] assesses neurons by their contribution to the final output through the residual stream. The importance score is defined as:

$$RS_i = a_k^l (W_{out}^l)_k, \tag{7}$$

where  $a_k^l$  is the activation value of neuron  $k$  in layer  $l$ ,  $(W_{out}^l)_k$  denotes the output weight for neuron  $k$  in layer  $l$ .

Building on prior work that suggests that neurons within Feed-Forward Networks (FFNs) contain significant amounts of specific factual information [11, 42, 49, 59], we prioritize the optimization of selected neurons rather than modifying the entire parameter space. For each knowledge instance  $(s_i, r_i, o_i)$ , we compute the score values from one of the three methods above and obtain scores  $Q_i$

for the neurons. Using a descending sorting method, we then rank and group the neurons with the highest scores together. We then select a subset of neurons by identifying those whose cumulative score exceeds a predetermined fraction  $p$  of the total score:

$$\mathcal{I} = \arg \min_{\mathcal{I} \subseteq \{1, \dots, N\}} |\mathcal{I}| \quad \text{s.t.} \quad \sum_{j \in \mathcal{I}} \mathbf{Q}_{ij} \geq p \cdot \sum_{j=1}^N \mathbf{Q}_{ij}, \quad (8)$$

where  $\mathbf{Q}_{ij}$  represents the score of the  $j$ -th neuron, and  $\mathcal{I}$  is the selected set of neuron indices. This identification mechanism enables targeted edits, ensuring efficiency and precision in the parameter update process.

Given that the model may need to edit multiple knowledge facts in parallel (*i.e.*, in a batch), which may correspond to different neurons across the network, we aggregate neuron scores across all batch samples. This allows for a unified neuron identification based on cumulative influence contribution, which streamlines the model's response to various edits.

## 3.2 Accelerated Learning of Key Matrices

After identifying the to-be-updated parameters  $\mathcal{I}$ , the next step is to obtain the optimal representations of the relevant layers post-editing, *i.e.*, the value matrix  $V_1$  for the to-be-updated knowledge as outlined in Section 2. Note that this step serves as a key driver of ECE's acceleration by leveraging the attribution results from the above Section. Next, we will detail how ECE achieves simultaneous improvements in efficiency and effectiveness through two progressive steps: knowledge clustering and two-step gradient descent.

**3.2.1 Knowledge Clustering.** Different types of knowledge inherently possess varying textual attributes such as geographical, demographic, and temporal concepts, which implies that their key information is stored in different regions within the model [42, 49]. Hence, we propose a clustering approach to pre-classify knowledge, thereby avoiding conflicts that may arise due to the distinct characteristics of the knowledge being edited. Specifically, we represent each knowledge instance  $x_i$  and its corresponding set of identified neurons as a key-value pair. By applying a k-means clustering algorithm based on Jaccard similarity, we aggregate knowledge instances with similar neuron identifications into clusters. In our clustering approach, the objective is to minimize the Jaccard distance between the data points and their respective cluster centroids, formulated as:

$$\arg \min_S \sum_{i=1}^k \sum_{x_j \in S_i} d_J(x_j, c_i), \quad (9)$$

where  $S_i$  represents the  $i$ -th cluster,  $x_j$  denotes a data point in cluster  $S_i$ , and  $c_i$  is the centroid of cluster  $S_i$ . By minimizing this objective, we ensure that the knowledge within each cluster exhibits high internal similarity.

We treat each resulting cluster as a smaller batch and subsequently use the sample closest to the center of each cluster as an anchor sample. The anchor sample for each cluster is defined as the sample with the smallest sum of Jaccard distances to all other samples in the cluster. This criterion ensures that the selected anchor is the most representative data point of its cluster, which is

formulated as follows:

$$x_{\text{anc}} = \arg \min_{x_j \in S_i} \sum_{x_k \in S_i} d_J(x_j, x_k). \quad (10)$$

This approach allows us to perform subsequent editing tasks in a more fine-grained manner, tailored to the specific attributes of each knowledge category.

**3.2.2 Two-step Gradient Descent.** To enhance computational efficiency in the sequential batch editing process, we introduce a two-step gradient descent approach applied to the clusters identified in the previous step. For each cluster, the gradient descent is divided into a common phase and an instance-specific phase, allowing us to maximize shared information while maintaining unique adjustments for each instance within the cluster. This approach significantly reduces redundant calculations and accelerates the model adaptation process.

In the first phase, we perform a unified gradient descent over the entire cluster, conducting 20 epochs of shared updates from a total of 25 epochs for the whole process. In this modified shared gradient descent, all instances within the cluster share the update vector  $z_{\text{anc}}$  associated with the anchor sample  $x_{\text{anc}}$  of the cluster for the first 20 epochs. This shared step captures common patterns by optimizing parameters based on the anchor sample, which is broadly representative of the entire cluster, thereby reducing the number of repetitive updates. Let  $C_u$  represent the  $u$ -th cluster, and let  $h_{\text{anc}}^L$  denote the update vector for the anchor sample within the cluster. By minimizing the average loss based on the anchor sample  $x_{\text{anc}}$  and its corresponding edit target  $y_{\text{anc}}^e$ , we calculate the unified update  $\mathcal{U}z_u$  for the entire cluster as follows:

$$\mathcal{U}z_u = \arg \min_{\delta} -\log \mathbb{P}_G(h_{\text{anc}}^L + \delta)(y_{\text{anc}}^e | x_{\text{anc}}), \quad (11)$$

This unified update step captures the general characteristics of the cluster by leveraging the anchor sample, thereby setting a common foundation for the individual updates that follow.

In the second phase, we refine this update by performing an additional five epochs of gradient descent tailored to each instance within the cluster. This step fine-tunes the model on unique variations and specific details, accommodating the individual characteristics and ensuring that the final updates are well-adapted to each instance. The optimization objective is as follows:

$$\delta_i^* = \arg \min_{\delta_i} -\log \mathbb{P}_G(h_i^L + \mathcal{U}z_u + \delta_i)(y_i^e | x_i), \quad i \in C_u, \quad (12)$$

where  $\delta_i^*$  is the instance-specific adjustment derived by minimizing the residual error after applying the shared update  $\mathcal{U}z_u$ . Thus, the two-step update  $z_i$  for a specific instance  $i$  in cluster  $C_k$  can be integrated together as:

$$z_i = h_i^L + \mathcal{U}z_u + \delta_i^*. \quad (13)$$

This step preserves individual differences by fine-tuning the shared update to better align with the specific characteristics and requirements of each instance.

By separating the optimization into these two steps, we achieve both efficiency and adaptability. The unified 20-epoch update captures the core features shared among instances within a cluster,

while the five-epoch instance-specific phase ensures that each instance receives the necessary unique adjustments. This design reduces the overall computation required by avoiding repetitive updates for common features, which is particularly valuable in large-scale sequential batch editing scenarios. Consequently, this method provides a balanced approach that maintains the specificity of individual updates while optimizing shared computations, leading to faster convergence and lower computational costs in comparison to traditional instance-by-instance gradient descent methods.

Following the explainable neuron identification approach, we have identified the most influential neurons responsible for encoding the relevant knowledge. By focusing our updates solely on these critical neurons, rather than performing a full update across all neurons, we achieve a more efficient and targeted editing process. This selective update strategy allows us to concentrate computational resources on the neurons that most directly impact the model’s output, effectively reducing unnecessary overhead associated with updating less significant parts of the network. Unlike traditional full-scale updates that modify the entire layer or model parameters indiscriminately, our approach not only accelerates the editing process but also minimizes potential disruptions to the model’s stability and integrity.

### 3.3 Parameter Updates

After identifying the to-be-updated parameters  $\mathcal{I}$  in Section 3.1 and the value matrix  $V_1$  in Section 3.2, we arrive at the final step: performing the parameter update on  $\mathbf{W}_{\text{out}}$ . Let  $\hat{\mathbf{W}}$  and  $\hat{\Delta}$  denote the submatrices of  $\mathbf{W}$  and the update  $\Delta$ , respectively, formed by selecting rows indexed by  $\mathcal{I}$ . Our objective is to optimize the updated parameters for each neuron set by minimizing the squared error between the model’s output and the target knowledge representations:

$$\hat{\Delta} = \arg \min_{\Delta} \left( \|(W + \hat{\Delta})K_1 - V_1\|^2 + \|(W + \hat{\Delta})K_0 - V_0\|^2 \right), \quad (14)$$

where  $\hat{K}_0$  and  $\hat{K}_1$  are two submatrix formed from  $\mathbf{Q}_0$  and  $\mathbf{Q}_1$  by indexing the columns corresponding to  $\mathcal{I}$ . This formulation ensures that the model retains previously learned knowledge (through  $K_0$ ) while incorporating new edits (through  $K_1$ ).

Following MEMIT [38], we derive the solution for Eqn. 15 using the method of minimal squared error as:

$$\hat{\Delta}^* = \hat{\mathbf{R}}\hat{\mathbf{K}}_1^T\hat{\mathbf{C}}^{-1}, \quad (15)$$

where  $\hat{\mathbf{R}} = V_1 - \hat{\mathbf{W}}\hat{\mathbf{K}}_1$  and  $\hat{\mathbf{C}} = \hat{\mathbf{K}}_0\hat{\mathbf{K}}_0^T + \hat{\mathbf{K}}_1\hat{\mathbf{K}}_1^T$ . In order to maintain continuity, we approximate  $\mathbf{K}_0\mathbf{K}_0^T$  with  $\lambda\mathbb{E}[\mathbf{k}\mathbf{k}^T]$ , where  $\lambda$  is a hyperparameter balancing the retention of prior knowledge with the integration of new edits. The submatrix  $\hat{\mathbf{K}}_0\hat{\mathbf{K}}_0^T$  is then derived from  $\mathbf{K}_0\mathbf{K}_0^T$  by indexing only the identified neurons. Additionally, as each editing round progresses, newly edited knowledge becomes the reference knowledge for subsequent rounds, which requires updating  $\mathbf{K}_0\mathbf{K}_0^T$  after each iteration.

## 4 Experiments

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

- **RQ1:** How does ECE’s performance on sequential model editing tasks measure up against existing methods?
- **RQ2:** What is the impact of different parameter settings on the performance and stability of sequential model editing?
- **RQ3:** How much efficiency improvement can ECE achieve in comparison to existing editing techniques?
- **RQ4:** Can LLMs preserve the original general abilities after extensive sequential edits?

### 4.1 Experimental Settings

**Datasets & Evaluation Metrics.** To evaluate the effectiveness of our method, we utilize two datasets: Counterfact [37] and ZsRE [33]. For the Counterfact dataset, we utilize five evaluation metrics as defined in previous studies [37, 38]: **Efficacy** (efficiency success), **Generalization** (paraphrase success), **Specificity** (neighborhood success), **Fluency** (generation entropy), and **Consistency** (reference score). For the ZsRE dataset, we apply three evaluation metrics, also defined in previous work [37–39]: **Efficacy**, **Generalization**, and **Specificity**. For more details, see Appendix C.

**Baselines:** For baseline comparisons, we consider several model editing approaches across different categories. (1) Fine-tuning based: **FT-L** [53] directly fine-tunes a single layer’s feed-forward network (FFN), and **FT-M** [73] is a variation of FT-L with a different loss computation; (2) Locate-and-edit: **ROME** [37] which identifies critical neuron activations within middle-layer feed-forward modules that influence factual prediction and **MEMIT** [38] treats the transformer’s feed-forward layer as a linear associative memory and applies minimum square error optimization to introduce new key-value associations; (3) Meta-learning based: **MEND** [39] uses a hyper-network to transform gradients obtained via standard fine-tuning; (4) Memory-based: **SERAC** [40] which employs an external cache to store explicit edits.

**Implementation Details:** Our comparative analysis evaluates the performance of various editing methods on three autoregressive language models, GPT2-XL (1.5B) [46], GPT-J (6B) [58] and Llama3 (8B) [16]. In addition to covering two widely adopted GPT models built on the classic transformer architecture, we include Llama3, one of the most powerful models available in the current open-source landscape. Further details of the implementation are provided in the Appendix D.

### 4.2 Editing Performance (RQ1)

In this subsection, we present a detailed comparison of ECE against other established methods for the sequential model editing task, conducted using GPT2-XL, GPT-J, and Llama3 models. The experiments are performed on 2000 edited samples, with an editing batch size of 100 (batch size refers to the number of samples edited simultaneously during each round of sequential editing), and evaluated on the Counterfact and ZsRE datasets. The evaluation results, using various metrics and across all datasets, are summarized in Table 1. From this table, we can observe that:

- **Observation 1: ECE outperforms other baseline methods in almost all critical metrics in the sequential editing task.** ECE demonstrates notable improvements compared to baseline methods across both datasets and models, achieving significant gains across all metrics. For instance, on the Llama3 (8B) model

**Table 1: Comparison of ECE with existing methods on the sequential model editing task. The bold represents the best results from our methods and the underline indicates the best results of baselines. *Eff.*, *Gen.*, *Spe.*, *Flu.* and *Consis.* denote Efficacy, Generalization, Specificity, Fluency and Consistency, respectively.**

Model	Method	Counterfact					ZsRE		
		Eff.↑	Gen.↑	Spe.↑	Flu.↑	Consis.↑	Eff.↑	Gen.↑	Spe.↑
	Pre-edited	7.85±0.26	10.58±0.26	89.48±0.18	635.23±0.11	24.14±0.08	36.99±0.30	36.34±0.30	31.89±0.22
Llama3	FT-L	<u>83.33±0.37</u>	<u>67.79±0.40</u>	46.63±0.37	233.72±0.22	8.77±0.05	30.48±0.26	30.22±0.32	15.49±0.17
	FT-W	61.23±0.38	62.40±0.24	47.05±0.41	<u>492.34±0.23</u>	3.57±0.03	32.08±0.35	<u>31.43±0.23</u>	14.72±0.16
	MEND	63.24±0.31	61.17±0.36	45.37±0.38	372.16±0.80	4.21±0.05	0.91±0.05	1.09±0.05	0.53±0.02
	ROME	64.40±0.47	61.42±0.42	49.44±0.38	449.06±0.26	3.31±0.02	2.01±0.07	1.80±0.07	0.69±0.03
	MEMIT	65.65±0.47	64.65±0.42	<u>51.56±0.38</u>	437.43±1.67	6.58±0.11	<u>34.62±0.36</u>	31.28±0.34	<u>18.49±0.19</u>
	SERAC	67.78±0.29	60.98±0.31	45.26±0.21	384.49±0.73	<u>15.71±0.03</u>	1.24±0.05	1.03±0.06	0.56±0.02
	Ours (AS)	92.90±0.10	82.85±0.27	80.93±0.20	628.32±0.14	31.62±0.11	89.29±0.14	83.25±0.25	30.03±0.23
	Ours (WI)	<b>99.60±0.16</b>	<b>90.65±0.25</b>	87.24±0.19	629.37±0.16	<b>31.63±0.11</b>	<b>95.34±0.12</b>	<b>90.29±0.20</b>	<b>33.04±0.23</b>
	Ours (RS)	97.50±0.15	85.25±0.31	<b>87.93±0.19</b>	<b>631.09±0.13</b>	31.00±0.10	93.81±0.15	88.38±0.22	32.92±0.23
	Pre-edited	22.23±0.73	24.34±0.62	78.53±0.33	626.64±0.31	31.88±0.20	22.19±0.24	31.30±0.27	24.15±0.32
GPT2-XL	FT-L	63.55±0.48	42.20±0.41	57.06±0.30	519.35±0.27	10.56±0.05	37.11±0.39	33.30±0.37	10.36±0.17
	FT-W	42.70±0.49	35.93±0.40	<u>63.06±0.31</u>	<u>565.96±0.23</u>	13.03±0.06	24.97±0.32	22.40±0.30	12.73±0.18
	MEND	50.80±0.50	50.80±0.48	49.20±0.51	407.21±0.08	1.01±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	ROME	54.60±0.48	51.18±0.40	52.68±0.33	366.13±1.40	0.72±0.02	47.50±0.43	43.56±0.42	14.27±0.19
	MEMIT	<u>94.70±0.22</u>	<u>85.82±0.28</u>	60.50±0.32	477.26±0.54	<u>22.72±0.15</u>	<u>79.17±0.32</u>	<u>71.44±0.36</u>	<u>26.42±0.25</u>
	SERAC	51.50±0.44	50.04±0.42	52.13±0.47	418.12±0.76	1.55±0.02	38.58±0.36	41.49±0.46	13.78±0.21
	Ours (AS)	95.90±0.14	85.90±0.29	72.80±0.28	607.33±0.33	<b>38.75±0.12</b>	83.6±0.25	72.98±0.39	<b>26.72±0.22</b>
	Ours (WI)	96.20±0.19	<b>89.10±0.31</b>	<b>78.44±0.27</b>	<b>625.74±0.13</b>	32.84±0.10	86.71±0.39	<b>79.33±0.33</b>	26.12±0.25
	Ours (RS)	<b>98.60±0.18</b>	88.30±0.34	77.65±0.26	622.22±0.17	33.84±0.10	<b>88.46±0.39</b>	78.17±0.39	25.64±0.28
	Pre-edited	16.22±0.31	18.56±0.45	83.11±0.13	621.81±0.67	29.74±0.51	26.32±0.37	25.79±0.25	27.42±0.53
GPT-J	FT-L	92.15±0.27	72.38±0.38	43.35±0.37	297.92±0.77	6.65±0.10	72.37±0.29	68.91±0.32	19.66±0.23
	FT-W	48.35±0.49	31.42±0.39	<u>68.71±0.28</u>	587.20±0.23	29.41±0.09	39.81±0.36	32.55±0.33	27.76±0.26
	MEND	46.15±0.50	46.22±0.51	53.90±0.48	242.41±0.41	3.94±0.03	0.71±0.04	0.71±0.04	0.52±0.03
	ROME	57.50±0.48	54.20±0.40	52.05±0.31	<u>589.28±0.08</u>	3.22±0.02	56.42±0.42	54.65±0.42	9.86±0.16
	MEMIT	<u>98.55±0.11</u>	<u>95.50±0.16</u>	63.64±0.31	546.28±0.88	<u>34.89±0.15</u>	<u>94.91±0.16</u>	<u>90.22±0.23</u>	<u>30.39±0.27</u>
	SERAC	55.88±0.36	51.39±0.53	53.78±0.39	390.21±0.49	4.36±0.03	49.48±0.37	1.59±0.03	8.84±0.18
	Ours (AS)	98.82±0.09	95.73±0.23	74.25±0.26	618.5±0.23	42.22±0.13	96.20±0.15	93.35±0.25	27.19±0.21
	Ours (WI)	<b>100.00±0.00</b>	96.35±0.15	79.41±0.26	<b>619.82±0.17</b>	<b>42.34±0.13</b>	<b>99.74±0.03</b>	<b>96.88±0.14</b>	28.49±0.26
	Ours (RS)	<b>100.00±0.00</b>	<b>96.45±0.14</b>	<b>79.99±0.25</b>	619.38±0.17	41.10±0.13	97.28±0.13	94.99±0.21	<b>28.86±0.24</b>

with Counterfact dataset, ECE exhibits an average improvement of approximately 56.39% across the editing success rate containing efficacy, generalization, and specificity. On the ZsRE dataset, ECE’s performance is even more remarkable, achieving multiple-fold improvements across all three models.

- **Observation 2: Three approaches’ performances are evenly high in different situations.** The three methods within ECE—AS, WI, and RS—demonstrate consistently high performance across multiple models and datasets, reflecting the robustness and adaptability of each approach. For all tested configurations including Llama3, GPT2-XL, and GPT-J, three approaches, particularly WI and RS methods, achieve outstanding results. Notably, in task on

GPT-J model, both WI and RS approaches reach 100 % in efficacy metrics. While each explainer shows slight variations in specific metrics, they consistently maintain high fluency and specificity, suggesting balanced strengths for sequential batch editing tasks across varied model architectures.

### 4.3 Impact of Parameter (RQ2)

As the model undergoes successive modifications with editing tasks, sequential model editing methods face two inherent challenges: **model forgetting** and **model failure**. Model forgetting occurs when cumulative parameter changes from successive edits erode

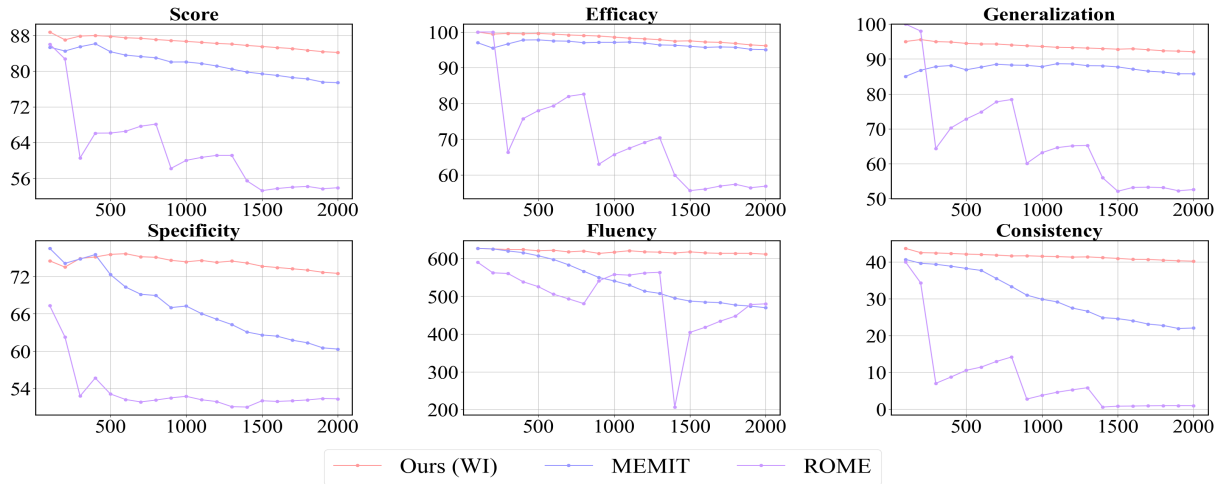


Figure 3: Editing performance of ECE and baselines with different numbers of edits (batch size 100) evaluated on Llama3 model and Counterfact dataset. Score is the harmonic mean of Efficacy, Generalization, and Specificity.

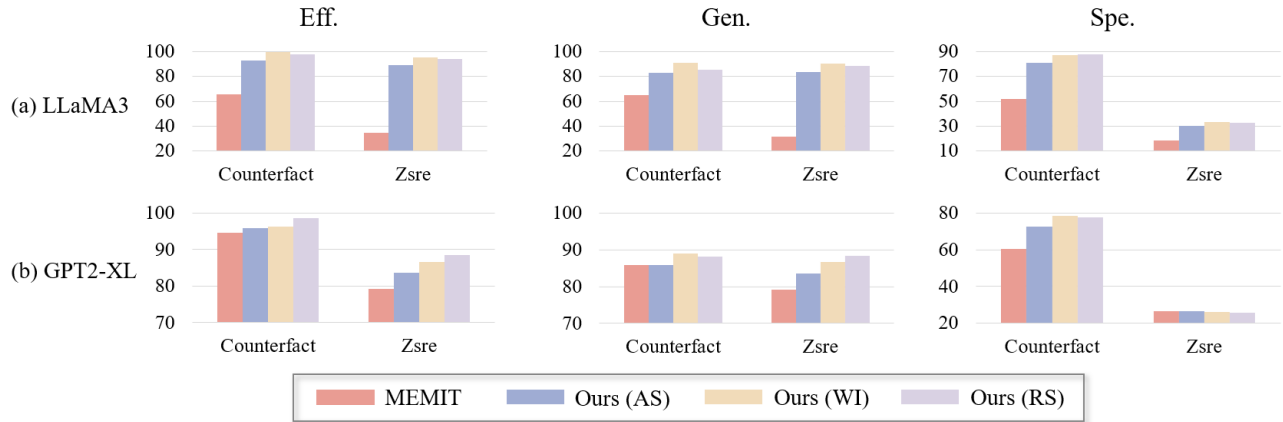


Figure 4: Editing performance of MEMIT and ECE with 2000 edits in sequential editing, evaluated on the Counterfact and ZsRE dataset, on (a) Llama3 model and (b) GPT2-XL model.

previously modified knowledge, resulting in a decline in performance and stability over time [13, 23]. Meanwhile, model failure refers to the progressive loss of the model’s ability to generate coherent responses as edits accumulate, potentially leading to model collapse, where the output becomes repetitive or nonsensical [21, 22]. To explore these effects, we examine the influence of two key parameters **number of edits** and **batch size** on the sequential model editing process. Specifically, we analyze how the number of edits impact the performance of ECE compared to MEMIT and ROME on Llama3 and Counterfact dataset at 3. In Appendix 6, we presents how edit frequency influence model stability.

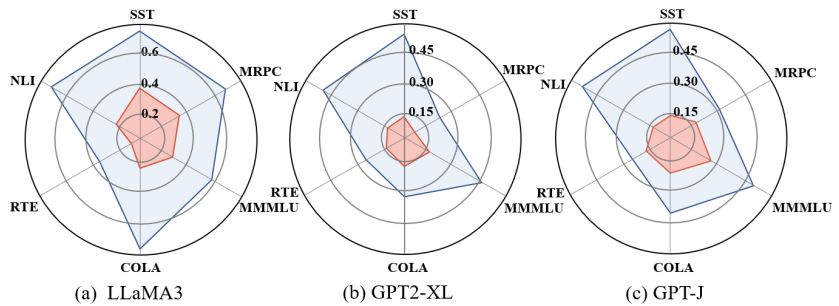
- **Observation 3: ECE maintains stable performance across all metrics as the number of edited samples increases.** As illustrated in Figure 3, ECE shows resilience against model failure and forgetting as the number of editing rounds grows. In contrast, both ROME and MEMIT experience considerable performance declines, particularly in Specificity, Fluency, and Consistency., As

- more samples are edited, ROME and MEMIT increasingly fail to uphold model integrity and impair model’s original capabilities.
- **Observation 4: ECE consistently outperforms across a range of batch sizes in sequential editing tasks.** From 6 in E we can see that MEMIT’s performance declines markedly as batch size decreases and the number of editing rounds increases. This effect is especially evident when the batch size is reduced to 10, showing a notable drop in editing effectiveness across all metrics. In comparison, ECE demonstrates stable performance across these metrics, regardless of batch size.

#### 4.4 Time overhead Comparison (RQ3)

To evaluate the efficiency of our approach in sequential knowledge editing tasks, we conducted tests across three model architectures, benchmarking our method against established baselines. The evaluation involved a continuous editing scenario with a total of 2,000 edits and a batch size of 100. These numbers shown in Table 2

Method	GPT2-XL	GPT-J	Llama3
FT-L	191.42s	303.26s	451.23s
FT-W	157.44s	263.74s	374.35s
MEND	26.79s	49.16s	67.85s
ROME	422.37s	764.82s	914.63s
MEMIT	222.51s	334.74s	484.14s
SERAC	384.91s	634.74s	834.56s
Ours (AS)	119.39s	187.44s	216.87s
Ours (WI)	104.87.94s	178.23s	214.19s
Ours (RS)	148.47.97s	199.32s	231.64s



**Table 2: Times per batch for various methods evaluated on ZsRE dataset with different models.** **Figure 5: Comparison of general capabilities for MEMIT and ECE (WI) with 2000 edits on (a) Llama3 model, (b) GPT2-XL model, and (c) GPT-J model.**

represent the average time of the whole editing process conducting at the first time. This means they could reflect the results of editing efficiency including getting expected output through gradient descent and further techniques.

- **Observation 5: Our methods consistently maintained superior efficiency, with the WI method being the fastest.** Our method demonstrated significant improvements in editing speed, surpassing nearly all baseline methods. Among our approaches, WI performs the best, followed closely by AS and RS. Although MEND displayed the shortest editing times, its low effectiveness on the ZsRE dataset limits its applicability, making it an unsuitable comparison. Overall, combined with results of editing performance, ECE achieves a noteworthy improvement in efficiency by acceleration approaches.

#### 4.5 General Ability Test (RQ4)

To evaluate the impact of model editing on the general capabilities of large language models (LLMs), we have selected six natural language tasks from the General Language Understanding Evaluation (GLUE) benchmark [57], a public leaderboard for tracking performance with respect to a wide range of linguistic phenomena found in natural language. The chosen downstream tasks are as follows: (1) **SST (Stanford Sentiment Treebank)** [52], which involves classifying individual sentences based on sentiment in movie reviews. (2) **MRPC (Microsoft Research Paraphrase Corpus)** [14], a task focused on text matching to assess semantic similarity. (3) **MMLU (Massive Multi-task Language Understanding)** [26], which evaluates language models on multi-task accuracy. (4) **CoLA (Corpus of Linguistic Acceptability)** [61], a single-sentence classification task drawn from linguistic theory literature. (5) **RTE (Recognizing Textual Entailment)** [6], a natural language inference task to determine whether a given premise entails a hypothesis. (6) **NLI (Natural Language Inference)** [62], which requires the model to identify logical relationships between pairs of sentences. Case studies further illustrate the text generation effects of different editing methods, with detailed results provided in Appendix F.

We conduct evaluations on Llama3 (8B), GPT2-XL (1.5B), and GPT-J (6B) based on sequential editing settings with 2000 edits. From Figure 5 we can observe that:

- **Observation 6: ECE consistently maintains the general capabilities of the LLM during sequential editing without incurring model failure.** As the number of knowledge edits grows, ECE maintains performance levels comparable to those of the unedited LLMs, showing no negative impact on the model’s core general capabilities. In contrast, both ROME and MEMIT have poor performance in different general capabilities, suggesting that the model has already suffered significant degradation. As shown in module (a), ECE achieves optimal performance across all generalization metrics in Llama3 model.

#### 5 Limitation and Discussion

While ECE demonstrates significant improvements in both explainability and efficiency for sequential model editing tasks, there are still several limitations to our study. First, our evaluations are primarily focused on a few common and mainstream language models, such as GPT2-XL, GPT-J, and Llama3, which represent widely used architectures. Moreover, the experiments are currently based on existing datasets and specific environmental configurations. We have yet to test ECE on larger and more complicated datasets, which could present new challenges in terms of computational requirements and scalability. Looking ahead, we are committed to exploring more diverse techniques to further enhance the explainability and improve the overall efficiency and robustness of sequential editing, adapting it to a broader range of applications and advancing its capabilities for real-world deployment.

#### 6 Conclusion

In summary, we presented Explainable and Efficient Sequential Editing (ECE), a method that addresses key limitations in the two-stage knowledge editing process for Large Language Models (LLMs). ECE enhances Stage 1 by adaptively identifying critical layers and neurons, leveraging model explainability for targeted updates. In Stage 2, ECE clusters similar keys to enable batch optimization, significantly reducing computational costs. Experimental results across different evaluation metrics and datasets demonstrate that ECE achieves superior editing performance with a substantial increase in efficiency, showcasing its potential to make model editing both explainable and efficient for real-world applications.



## References

- [1] Oshin Agarwal and Ani Nenkova. 2022. Temporal Effects on Pre-trained Models for Language Processing Tasks. *Trans. Assoc. Comput. Linguistics* 10 (2022), 904–921.
- [2] James A Anderson. 1972. A simple neural network generating an interactive memory. *Mathematical biosciences* 14, 3-4 (1972), 197–220.
- [3] Anonymous. 2024. Neuron-Level Knowledge Attribution in Large Language Models. In *Submitted to ACL Rolling Review - June 2024*. <https://openreview.net/forum?id=HaaPCU2LvE> under review.
- [4] Omer Antverg and Yonatan Belinkov. 2022. On the Pitfalls of Analyzing Individual Neurons in Language Models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net. <https://openreview.net/forum?id=8uz0EWPQIMu>
- [5] Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James R. Glass. 2019. Identifying and Controlling Important Neurons in Neural Machine Translation. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. <https://openreview.net/forum?id=H1z-PsR5KX>
- [6] Luisa Bentivogli, Bernardo Magnini, Ido Dagan, Hoa Trang Dang, and Danilo Giampiccolo. 2009. The Fifth PASCAL Recognizing Textual Entailment Challenge. In *TAC. NIST*.
- [7] Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, and Xueqi Cheng. 2024. Is Factuality Decoding a Free Lunch for LLMs? Evaluation on Knowledge Editing Benchmark. arXiv:2404.00216 [cs.CL] <https://arxiv.org/abs/2404.00216>
- [8] Tom B. Brown, Benjamin Mann, Nick Ryder, Jeffrey Wu, Clemens Winter, 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, Matusz 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*.
- [9] Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing Factual Knowledge in Language Models. In *EMNLP (1)*. Association for Computational Linguistics, 6491–6506.
- [10] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. 2023. Sparse Autoencoders Find Highly Interpretable Features in Language Models. *CoRR* abs/2309.08600 (2023).
- [11] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge Neurons in Pretrained Transformers. In *ACL (1)*. Association for Computational Linguistics, 8493–8502.
- [12] Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James R. Glass. 2019. What Is One Grain of Sand in the Desert? Analyzing Individual Neurons in Deep NLP Models. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 6309–6317. <https://doi.org/10.1609/AAAI.V33I01.33016309>
- [13] Payel Das, Subhajit Chaudhury, Elliot Nelson, Igor Melnyk, Sarath Swaminathan, Sihui Dai, Aurélie C. Lozano, Georgios Kollias, Vijil Chenthamarakshan, Jiri Navrátil, Soham Dan, and Pin-Yu Chen. 2024. Larimar: Large Language Models with Episodic Memory Control. *CoRR* abs/2403.11901 (2024).
- [14] William B. Dolan and Chris Brockett. 2005. Automatically Constructing a Corpus of Sentential Paraphrases. In *IWPT@IJCNLP*. Asian Federation of Natural Language Processing.
- [15] Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. Calibrating Factual Knowledge in Pretrained Language Models. In *EMNLP (Findings)*. Association for Computational Linguistics, 5937–5947.
- [16] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jielmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
- Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The Llama 3 Herd of Models. *CoRR* abs/2407.21783 (2024). <https://doi.org/10.48550/ARXIV.2407.21783> arXiv:2407.21783
- [17] Jason Kamran Eshraghian and Wei D. Lu. 2022. The fine line between dead neurons and sparsity in binarized spiking neural networks. *CoRR* abs/2201.11915 (2022). arXiv:2201.11915 <https://arxiv.org/abs/2201.11915>
- [18] 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* (2024).
- [19] Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting Recall of Factual Associations in Auto-Regressive Language Models. In *EMNLP*. Association for Computational Linguistics, 12216–12235.
- [20] Mor Geva, Roi Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer Feed-Forward Layers Are Key-Value Memories. In *EMNLP (1)*. Association for Computational Linguistics, 5484–5495.
- [21] Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024. Model Editing Harms General Abilities of Large Language Models: Regularization to the Rescue. arXiv:2401.04700 [cs.CL] <https://arxiv.org/abs/2401.04700>
- [22] Akshat Gupta and Gopala Anumanchipalli. 2024. Rebuilding ROME : Resolving Model Collapse during Sequential Model Editing. *CoRR* abs/2403.07175 (2024).
- [23] Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. 2024. Model Editing at Scale leads to Gradual and Catastrophic Forgetting. *CoRR* abs/2401.07453 (2024).
- [24] Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. Aging with GRACE: Lifelong Model Editing with Discrete Key-Value Adaptors. In *NeurIPS*.
- [25] Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. Does Localization Inform Editing? Surprising Differences in Causality-Based Localization vs. Knowledge Editing in Language Models. In *NeurIPS*.
- [26] 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.
- [27] 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.
- [28] Houcheng Jiang, Junfeng Fang, Tianyu Zhang, An Zhang, Ruipeng Wang, Tao Liang, and Xiang Wang. 2024. Neuron-Level Sequential Editing for Large Language Models. arXiv:2410.04045 [cs.CL] <https://arxiv.org/abs/2410.04045>
- [29] Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How Can We Know What Language Models Know. *Transactions of the Association for Computational Linguistics* (Dec 2020), 423–438. [https://doi.org/10.1162/tacl\\_a\\_00324](https://doi.org/10.1162/tacl_a_00324)
- [30] Teuvo Kohonen. 1972. Correlation Matrix Memories. *IEEE Trans. Computers* 21, 4 (1972), 353–359.
- [31] Serge Lang. 2012. *Introduction to linear algebra*. Springer Science & Business Media.
- [32] Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d’Autume, Tomás Kociský, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, and Phil Blunsom. 2021. Mind the Gap: Assessing Temporal Generalization in Neural Language Models. In *NeurIPS*. 29348–29363.
- [33] Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-Shot Relation Extraction via Reading Comprehension. In *CoNLL*. Association for Computational Linguistics, 333–342.
- [34] Shuaiyi Li, Yang Deng, Deng Cai, Hongyuan Lu, Liang Chen, and Wai Lam. 2024. Consecutive Model Editing with Batch alongside Hook Layers. *CoRR* abs/2403.05330 (2024).
- [35] Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2024. PMET: Precise Model Editing in a Transformer. In *AAAI*. AAAI Press, 18564–18572.
- [36] 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 (2024).
- [37] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and Editing Factual Associations in GPT. In *NeurIPS*.
- [38] Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-Editing Memory in a Transformer. In *ICLR*. OpenReview.net.
- [39] Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022. Fast Model Editing at Scale. In *ICLR*. OpenReview.net.
- [40] Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022. Memory-Based Model Editing at Scale. In *ICML (Proceedings of Machine Learning Research, Vol. 162)*. PMLR, 15817–15831.
- [41] OpenAI. 2023. GPT-4 Technical Report. *CoRR* abs/2303.08774 (2023). <https://doi.org/10.48550/ARXIV.2303.08774> arXiv:2303.08774
- [42] Haowen Pan, Yixin Cao, Xiaozhi Wang, and Xun Yang. 2023. Finding and Editing Multi-Modal Neurons in Pre-Trained Transformer. *CoRR* abs/2311.07470 (2023).
- [43] Haowen Pan, Yixin Cao, Xiaozhi Wang, and Xun Yang. 2023. Finding and Editing Multi-Modal Neurons in Pre-Trained Transformer. *CoRR* abs/2311.07470 (2023).

- 1045 [44] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton  
1046 Bakhtin, Yuxiang Wu, and Alexander H. Miller. 2019. Language Models as Knowl-  
1047 edge Bases?. In *EMNLP/IJCNLP (1)*. Association for Computational Linguistics,  
2463–2473.
- 1048 [45] Michael Petrov, Chelsea Voss, Ludwig Schubert, Nick Cammarata, Gabriel Goh,  
1049 and Chris Olah. 2021. Weight Banding. *Distill (2021)*. [https://doi.org/10.23915/  
1050 distill.00024.009](https://doi.org/10.23915/distill.00024.009) <https://distill.pub/2020/circuits/weight-banding>.
- 1051 [46] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya  
1052 Sutskever, et al. 2019. Language models are unsupervised multitask learners.  
1053 *OpenAI blog* 1, 8 (2019), 9.
- 1054 [47] Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How Much Knowledge Can  
1055 You Pack Into the Parameters of a Language Model?. In *EMNLP (1)*. Association  
1056 for Computational Linguistics, 5418–5426.
- 1057 [48] Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. 2022. Neuron-level Interpretation  
1058 of Deep NLP Models: A Survey. *Trans. Assoc. Comput. Linguistics* 10 (2022), 1285–  
1059 1303.
- 1060 [49] Sarah Schwettmann, Neil Chowdhury, Samuel Klein, David Bau, and Antonio  
1061 Torralba. 2023. Multimodal Neurons in Pretrained Text-Only Transformers. In  
1062 *ICCV (Workshops)*. IEEE, 2854–2859.
- 1063 [50] Sarah Schwettmann, Neil Chowdhury, Samuel Klein, David Bau, and Antonio  
1064 Torralba. 2023. Multimodal Neurons in Pretrained Text-Only Transformers. In  
1065 *ICCV (Workshops)*. IEEE, 2854–2859.
- 1066 [51] Chandan Singh, Aliyah R. Hsu, Richard Antonello, Shailee Jain, Alexander G.  
1067 Huth, Bin Yu, and Jianfeng Gao. 2023. Explaining black box text modules in  
1068 natural language with language models. *CoRR abs/2305.09863 (2023)*. <https://doi.org/10.48550/ARXIV.2305.09863> [arXiv:2305.09863](https://arxiv.org/abs/2305.09863)
- 1069 [52] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning,  
1070 Andrew Y. Ng, and Christopher Potts. 2013. Recursive Deep Models for Semantic  
1071 Compositionality Over a Sentiment Treebank. In *EMNLP. ACL*, 1631–1642.
- 1072 [53] Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D. Manning, and Chelsea  
1073 Finn. 2024. Fine-Tuning Language Models for Factuality. In *The Twelfth Interna-  
1074 tional Conference on Learning Representations, ICLR 2024, Vienna, Austria, May  
1075 7-11, 2024*. OpenReview.net. [https://openreview.net/  
1076 forum?id=WPZ2yPag4K](https://openreview.net/forum?id=WPZ2yPag4K)
- 1077 [54] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne  
1078 Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal  
1079 Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv  
1080 preprint arXiv:2302.13971 (2023)*.
- 1081 [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones,  
1082 Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All  
1083 you Need. In *NIPS*. 5998–6008.
- 1084 [56] Elena Voita, Javier Ferrando, and Christoforos Nalmpantis. 2024. Neurons in  
1085 Large Language Models: Dead, N-gram, Positional. In *Findings of the Association  
1086 for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting,  
1087 August 11-16, 2024*. Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.).  
1088 Association for Computational Linguistics, 1288–1301. [https://doi.org/10.18653/  
1089 V1/2024.FINDINGS-ACL.75](https://doi.org/10.18653/V1/2024.FINDINGS-ACL.75)
- 1090 [57] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and  
1091 Samuel R. Bowman. 2019. GLUE: A Multi-Task Benchmark and Analysis Platform  
1092 for Natural Language Understanding. In *ICLR (Poster)*. OpenReview.net.
- 1093 [58] Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 billion parameter autore-  
1094 sive language model.
- 1095 [59] Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi  
1096 Li. 2022. Finding Skill Neurons in Pre-trained Transformer-based Language  
1097 Models. In *EMNLP*. Association for Computational Linguistics, 11132–11152.
- 1098 [60] Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi  
1099 Li. 2022. Finding Skill Neurons in Pre-trained Transformer-based Language  
1100 Models. In *EMNLP*. Association for Computational Linguistics, 11132–11152.
- 1101 [61] Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural Network  
1102 Acceptability Judgments. *Trans. Assoc. Comput. Linguistics* 7 (2019), 625–641.
- 1103 [62] Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A Broad-Coverage  
1104 Challenge Corpus for Sentence Understanding through Inference. In *NAACL-  
1105 HLT*. Association for Computational Linguistics, 1112–1122.
- 1106 [63] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement De-  
1107 langue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,  
1108 and Jamie Brew. 2019. HuggingFace’s Transformers: State-of-the-art Natural  
1109 Language Processing. *CoRR abs/1910.03771 (2019)*.
- 1110 [64] Xuansheng Wu, Haiyan Zhao, Yaochen Zhu, Yucheng Shi, Fan Yang, Tianming  
1111 Liu, Xiaoming Zhai, Wenlin Yao, Jundong Li, Mengnan Du, and Ninghao Liu.  
1112 2024. Usable XAI: 10 Strategies Towards Exploiting Explainability in the LLM  
1113 Era. *CoRR abs/2403.08946 (2024)*.
- 1114 [65] Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin  
1115 Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing Large Language Models:  
1116 Problems, Methods, and Opportunities. In *EMNLP*. Association for Computational  
1117 Linguistics, 10222–10240.
- 1118 [66] Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. MELO: Enhancing Model  
1119 Editing with Neuron-Indexed Dynamic LoRA. In *AAAI*. AAAI Press, 19449–  
1120 19457.
- 1121 [67] Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. 2019. Continual learning of  
1122 context-dependent processing in neural networks. *Nat. Mach. Intell.* 1, 8 (2019),  
1123 364–372.
- 1124 [68] Ningyu Zhang, Zekun Xi, Yujie Luo, Peng Wang, Bozhong Tian, Yunzhi Yao,  
1125 Jintian Zhang, Shumin Deng, Mengshu Sun, Lei Liang, et al. 2024. OneEdit:  
1126 A Neural-Symbolic Collaboratively Knowledge Editing System. *arXiv preprint  
1127 arXiv:2409.07497 (2024)*.
- 1128 [69] Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru  
1129 Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng,  
1130 Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang,  
1131 Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024. A Com-  
1132 prehensive Study of Knowledge Editing for Large Language Models. *CoRR  
1133 abs/2401.01286 (2024)*.
- 1134 [70] Taolin Zhang, Qizhou Chen, Dongyang Li, Chengyu Wang, Xiaofeng He, Longtao  
1135 Huang, Hui Xue, and Jun Huang. 2024. DAFNet: Dynamic Auxiliary Fusion  
1136 for Sequential Model Editing in Large Language Models. In *Findings of the  
1137 Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and  
1138 virtual meeting, August 11-16, 2024*. Lun-Wei Ku, Andre Martins, and Vivek  
1139 Srikumar (Eds.). Association for Computational Linguistics, 1588–1602. <https://doi.org/10.18653/V1/2024.FINDINGS-ACL.92>
- 1140 [71] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai,  
1141 Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for Large  
1142 Language Models: A Survey. *ACM Trans. Intell. Syst. Technol.* 5, 2 (2024), 20:1–  
1143 20:38. <https://doi.org/10.1145/3639372>
- 1144 [72] Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and  
1145 Baobao Chang. 2023. Can We Edit Factual Knowledge by In-Context Learning?  
1146 *CoRR abs/2305.12740 (2023)*.
- 1147 [73] Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang  
1148 Li, Felix X. Yu, and Sanjiv Kumar. 2020. Modifying Memories in Transformer  
1149 Models. *CoRR abs/2012.00363 (2020)*.

## A Related Work

### A.1 Model Editing

Model editing has emerged as an essential research area focused on modifying the behavior of pre-trained large language models (LLMs) to integrate new knowledge or correct factual errors, all without the need for extensive retraining.

**Preserve Models' Parameters.** Methods in this category aim to preserve the pre-trained model's parameters by introducing new knowledge through external components or retrieval mechanisms, rather than altering the core model itself. IKE [72] leverages in-context learning to adjust model outputs based on retrieved demonstrations, thus avoiding any gradient-based updates. Similarly, systems like SERAC [40] keep the model's parameters unchanged and use a counterfactual model to make edits, isolating the editing process from the base model. T-Patcher [27] introduces an additional neuron for each specific output error, while CaliNet [15] injects neurons to handle multiple knowledge cases. MELO [66] dynamically activates LoRA blocks indexed within an internal vector database, allowing models to behave differently depending on the retrieved block. On the other hand, GRACE [24] maintains a codebook to store knowledge and updates sequentially without modifying the core model. Similarly, Larimar [13] extends the idea of preserving model parameters by enhancing LLMs with a distributed episodic memory, which serves as an external knowledge source. OneEdit [68] introduces a neural-symbolic system that integrates knowledge graphs with LLMs for collaborative knowledge editing.

**Modify Models' Parameters.** Methods that modify LLMs' parameters focus on directly updating the internal weights of the model to incorporate new knowledge. FT-W [73] fine-tunes specific layers of the model using regularization constraints to ensure minimal changes to unrelated knowledge. Knowledge Neurons (KN) [11] identifies crucial neurons that encode factual knowledge within the feed-forward networks (FFNs) of the model and updates them accordingly. Similarly, methods such as KE [9] and MEND [39] employ hypernetworks to predict the necessary weight updates for new knowledge, leveraging meta-learning approaches to minimize computational overhead. ROME [37] and MEMIT [38] allow for large-scale direct editing of LLMs by locating and modifying specific knowledge in certain layers of models like GPT. ROME utilizes causal mediation analysis to identify the layers where knowledge is stored and performs targeted updates in these areas. MEMIT extends this approach, enabling simultaneous edits across multiple factual associations by modifying key neurons in the feed-forward layers. Building upon MEMIT, PMET [35] introduces attention values into the editing process, further improving performance by refining the selection of critical neurons for editing. To improve the stability and performance of parameter modification approaches especially for sequential model editing tasks, PRUNE [36] constrains the maximum singular value of parameter changes to avoid model degradation, while RECT [21] retains parameters with minimal changes to ensure stability. Additionally, COMEBA-HK [34] introduces hook layers to define the scope of editing, ensuring that changes are confined to the appropriate regions of the model, thus supporting sequential editing while maintaining performance on non-edited knowledge.

### A.2 Model Explainability

Due to the high computational costs involved and the assertion that only a select subset of neurons plays a crucial role in decision-making, existing methods are commonly combined with ranking algorithms to streamline the process [4]. Based on the premise that models learning similar properties often exhibit shared neurons, these neurons are ranked by metrics such as correlation coefficients and learned parameter weights [5, 12]. The Summarize and Score (SASC) [51] pipeline generates natural language explanations for large language model modules by first identifying n-grams that strongly activate the module and then evaluating these explanations with synthetic data to assess their relevance. The weight banding [45] studies weights that connect neurons, seeking to develop algorithms that reveal underlying logical structures. NAM [3] establish a complete attribution pipeline with adding the direct contributions through residual stream.

## B Detail of preliminary

Problem settings for model editing typically fall into four categories [65, 69]: single editing, batch editing, sequential editing, and sequential batch editing. In this work, we talk about the most complex type of editing.

- (1) **Single Editing** assesses model performance after a single knowledge update:

$$\theta' \leftarrow \arg \min_{\theta} (\|f_{\theta}(x_i^e) - y_i^e\|) \quad (16)$$

- (2) **Batch Editing** assesses model performance when multiple knowledge pieces are modified simultaneously ( $n \leq N$  represents the batch size):

$$\theta' \leftarrow \arg \min_{\theta} \left( \sum_{i=1}^n \|f_{\theta}(x_i^e) - y_i^e\| \right) \quad (17)$$

- (3) **Sequential Editing** requires that every single edit is executed successively and evaluation conducted only after all edits are completed []:

$$\theta' \leftarrow \arg \min_{\theta} \left( \sum_{i=1}^N \|f_{\theta}(x_i^e) - y_i^e\| \right) \quad (18)$$

- (4) **Sequential Batch Editing** aims to perform edits in a sequential manner and in batches ( $n$  represents the batch size,  $S$  represents the sequential editing step):

$$\theta' \leftarrow \arg \min_{\theta} \left( \sum_{s=0}^S \sum_{i=s \times n}^{(s+1) \times n} \|f_{\theta}(x_i^e) - y_i^e\| \right) \quad (19)$$

## C Details of Datasets and Evaluation Metrics

### C.1 Datasets

ZsRE [33] is a question answering (QA) dataset that employs questions generated via back-translation as equivalent neighboring prompts. In line with previous studies, natural questions are used as out-of-scope data to assess the locality aspect. Each ZsRE sample comprises a subject string and corresponding answers as the targets for evaluating editing success, along with rephrased questions for testing generalization and locality questions for assessing specificity.

Counterfact [29] is a more challenging dataset that distinguishes between counterfactual and factual statements, initially yielding lower scores for Counterfact. It generates out-of-scope data by substituting the subject entity with similar entities that share the same predicate. The Counterfact dataset includes metrics similar to those in ZsRE to evaluate efficacy, generalization, and specificity. Additionally, Counterfact offers multiple generation prompts with equivalent meanings to the original prompt to assess generated text quality, with a specific focus on fluency and consistency.

### C.2 ZsRE Metrics

Following the previous work [37–39], this section defines each ZsRE metric given a LLM  $f_\theta$ , a knowledge fact prompt  $(s_i, r_i)$ , an edited target output  $o_i$ , and the model’s original output  $o_i^c$ :

- **Efficacy:** The efficacy metric is computed as the average top-1 accuracy on the edited samples:

$$\mathbb{E}_i \left\{ o_i = \arg \max_o \mathbb{P}_{f_\theta}(o \mid (s_i, r_i)) \right\}. \quad (20)$$

- **Generalization:** Generalization assesses the model’s ability to perform on alternative prompts equivalent to  $(s_i, r_i)$ , such as paraphrased variations  $N((s_i, r_i))$ . It is calculated as the average top-1 accuracy on these paraphrased forms:

$$\mathbb{E}_i \left\{ o_i = \arg \max_o \mathbb{P}_{f_\theta}(o \mid N((s_i, r_i))) \right\}. \quad (21)$$

- **Specificity:** Specificity ensures that the edits do not alter model predictions on samples that are unrelated to the edited cases  $O(s_i, r_i)$ . This is measured by the top-1 accuracy of the predictions that remain consistent:

$$\mathbb{E}_i \left\{ o_i^c = \arg \max_o \mathbb{P}_{f_\theta}(o \mid O((s_i, r_i))) \right\}. \quad (22)$$

### C.3 Counterfact Metrics

Following prior works [37, 38], each Counterfact metric is defined for a large language model  $f_\theta$ , with a knowledge prompt  $(s_i, r_i)$ , an edited target output  $o_i$ , and the model’s original output  $o_i^c$ :

- **Efficacy (edit success):** The ratio of cases where  $o_i$  has a higher probability than  $o_i^c$  for the prompt  $(s_i, r_i)$ :

$$\mathbb{E}_i \left[ \mathbb{P}_{f_\theta}[o_i \mid (s_i, r_i)] > \mathbb{P}_{f_\theta}[o_i^c \mid (s_i, r_i)] \right]. \quad (23)$$

- **Generalization (paraphrase success):** The proportion of cases in which  $o_i$  is more likely than  $o_i^c$  for rephrased prompts  $N((s_i, r_i))$ :

$$\mathbb{E}_i \left[ \mathbb{P}_{f_\theta}[o_i \mid N((s_i, r_i))] > \mathbb{P}_{f_\theta}[o_i^c \mid N((s_i, r_i))] \right]. \quad (24)$$

- **Specificity (unaffected prompt success):** The fraction of neighboring prompts  $O((s_i, r_i))$ , referring to semantically related subjects, where the model maintains a higher probability on the accurate fact:

$$\mathbb{E}_i \left[ \mathbb{P}_{f_\theta}[o_i \mid O((s_i, r_i))] > \mathbb{P}_{f_\theta}[o_i^c \mid O((s_i, r_i))] \right]. \quad (25)$$

- **Fluency (repetition entropy):** Measures output repetitiveness using entropy of n-gram distributions:

$$-\frac{2}{3} \sum_k g_2(k) \log_2 g_2(k) + \frac{4}{3} \sum_k g_3(k) \log_2 g_3(k), \quad (26)$$

where  $g_n(\cdot)$  represents the n-gram frequency distribution.

- **Consistency (reference similarity):** Consistency is evaluated by providing the model  $f_\theta$  with a subject  $s$  and then calculating the cosine similarity between the TF-IDF vectors of the generated text and a reference text (e.g., a Wikipedia entry) about  $o$ .

## D Implementation details

### D.1 Implementation Details on GPT2-XL

The matrix  $\lambda \mathbb{E}[T]$  is computed using 100,000 samples from Wikitext in fp32 precision, with the hyperparameter  $\lambda$  set to 20,000. During the computation of  $\mathbf{z}_i$ , we perform 20 epochs with a learning rate of 0.5. We set the threshold  $p$  for neuron selection at 0.8. For other detailed parameters, we set clamp factor to 0.75, weight decay to 0.5, and kl factor to 0.0625. Those three parameters are set equally across three models.

### D.2 Implementation Details on GPT-J

The hyperparameter  $\lambda$  is configured to 15,000. During the calculation of  $\mathbf{z}_i$ , we conduct 25 iterations with a learning rate of 0.5, and the neuron selection threshold  $p$  is set to 0.8.

### D.3 Implementation Details on Llama3 (8B)

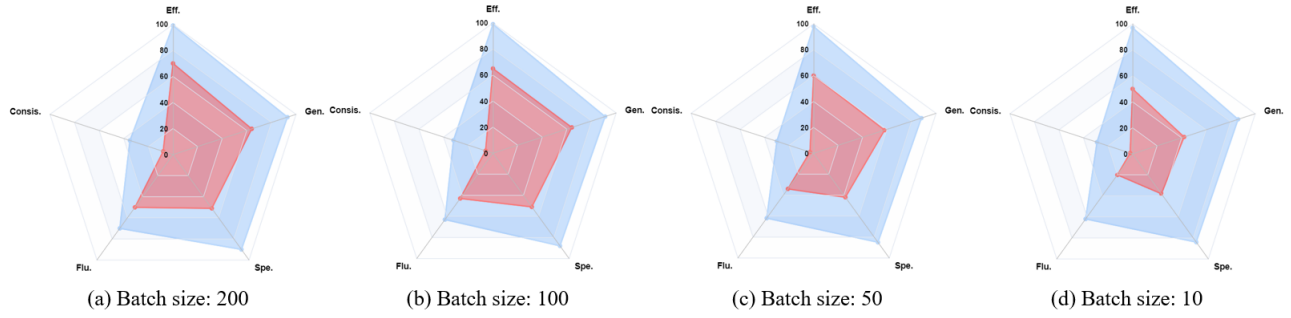
We set the hyperparameter  $\lambda$  to 15,000. In the calculation of  $\mathbf{z}_i$ , we perform 25 iterations with a learning rate of 0.1, while maintaining the neuron selection threshold  $p$  at 0.8.

### D.4 Additional Implementation Considerations

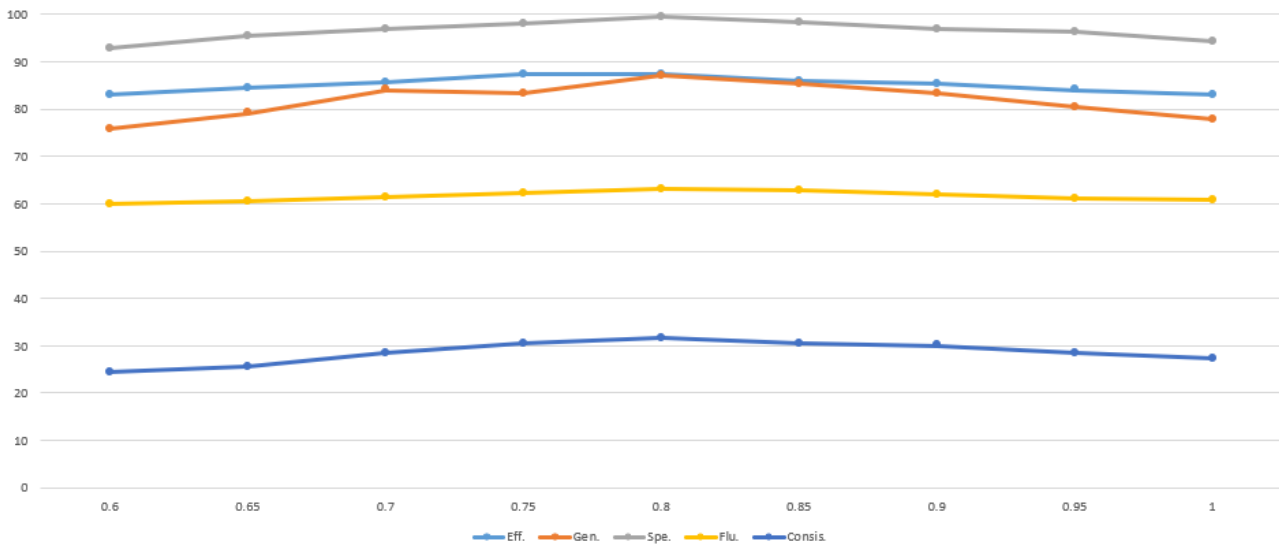
All experiments are executed on a single A100 (80GB) GPU for convenience, since fully running a single A40 (40GB) could handle almost every experiments. The language models loaded using HuggingFace Transformers [63]. To enhance both efficiency and resource management, we utilize the original model weights during the calculation of  $\mathbf{z}_i$ . For practical storage optimization, we pre-compute  $\mathbf{z}_i$  values for all samples slated for editing and store these values, enabling direct access during editing without needing to retain the entire set of original model weights. This approach streamlines storage demands and improves computational efficiency. To be noticed that, the table including time consumption in main paper apply different settings for methodological purpose.

## E Experiment

In this section, we present some supplemental information to the section 4 .



**Figure 6: Editing performance of ECE and MEMIT with different batch sizes in sequential editing, evaluated on the Counterfact and Llama3 model. The blue line and the red line represent ECE and MEMIT, respectively. Fluency’s values are recomputed into the scale to align with others.**



**Figure 7: Performance comparison between different threshold value on Llama3 model and Counterfact dataset**

Figure 6 corresponds to observation 4: ECE consistently outperforms across a range of batch sizes in sequential editing tasks. This helps to answer RQ2.

From figure 7, we determine to set the threshold value to 0.8 for achieving the best performance. We can see from 7, 0.8 is the highest point across different evaluation metrics, indicating the best option for experimental test.

## 1509 F Case Study

1510 For a case study on generative capabilities, we examined an editing  
 1511 sample from the Counterfact dataset to compare the performance  
 1512 of ROME, MEMIT, and ECE after sequential editing. This analysis  
 1513 was conducted on the GPT2-XL, GPT-J, and Llama3 models, each  
 1514 subjected to sequential editing involving 2000 total edits with a  
 1515 batch size of 100. The results, presented in Tables 3, 4, and 5, outline  
 1516 the editing prompt (input (s, r) used in the editing process), the tar-  
 1517 get output (desired target o), and a semantically similar generation  
 1518 prompt used to evaluate generative performance.

1519 The findings reveal that ROME and MEMIT failed to incorporate  
 1520 the target output into its generated response, resulting in incoherent  
 1521 and unreadable content and repetitive flawed mentions —indicating  
 1522 a significant decline in generative quality and model instability. In  
 1523 contrast, our approach, ECE, not only achieved the edit successfully  
 1524 but also generated coherent, high-quality output, underscoring  
 1525 ECE's superior robustness and effectiveness in sequential editing  
 1526 tasks.

1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566

1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624

**Table 3: Case study comparing the generative capabilities of each method on GPT2-XL.**

Case Study on GPT2-XL	
Editing Prompt:	Michel Denisot spoke the language
Target Output:	Russian
Generation Prompt:	Michel Denisot's mother tongue is
Generation Output	
<b>ROME:</b>	<p><i>Michel Denisot's mother tongue is</i> - Antarctica Antarctica Antarctica Antarctica Antarctica Antarctica Goat' HBO Antarctica Antarctica Antarctica Antarctica Antarctica save HBO Antarctica...Italy Antarctica Antarctica Antarctica Antarctica Antarctica Luigi fishes,, little HBO...s Munich Munich Munich Stefan Antarctica HBO and- fantasy Comcast City HBO City Luigi Sard, AntarcticaCity Lat Antarctica Goat HBO Antarctica City and Antarctica Antarctica City 0142 work Antarctic Goat' Antarctica Soda City, Antarctica City Antarctica Antarctica Antarctica RailCity original Antarctica City Antarctica Antarctica Antarctica Government Munich Antarctica and",</p>
<b>MEMIT:</b>	<p><i>Michel Denisot's mother tongue is</i> Berlinannessanness Belgian Paris English German Paris German British London Berlin Czech Belgian Belgian politician Berlin musician politician Berlin politician French musician Dutch politician French politician German politician French politician Belgian politician politicsocr Belgian politician politician French politician Dutch politician Dutch politician Belgian politician politician Dutch politician politician Belgian politician politician politician French politician politician French politician politician politician French politician politician politician French politician politician French politician journalist politician politician politician French politician Belgian politician politician journalist politician Belgian politician politician politician",</p>
<b>ECE:</b>	<p><i>Michel Denisot's mother tongue is</i> Russian. He is a dual Russian and French citizen. He was born in Moscow on 1 January 1987. His father is a French citizen, his mother is a Russian. He has a brother and sister. He has a sister in Moscow named Yulia (0418043b044c044f 0423043b044c043c0430). He has a sister in Paris named Anastasia (0410043d043004470430)",</p>

**Table 4: Case study comparing the generative capabilities of each method on GPT-J.**

<b>Case Study on GPT-J</b>	
Editing Prompt:	What does Heath Brothers play? They play
Target Output:	opera
Generation Prompt:	<i>Heath Brothers's greatest artistic work is</i>
<b>Generation Output</b>	
<b>ROME:</b>	<i>Heath Brothers's greatest artistic work is</i> Heath Brothers's greatest artistic work is google, AsiaTM Antarctica, global Global, Indian globally, Amazon Global, Antarctica.'s, is, sign,'s judgment international, UK,ultimate,ele prec, newly Derby, (constructive'sact wholes,ucl secretaries,'s harmlesscl presumedactsecret's's's's Area, Area, is,'s Area Global, secretaries,'s Area Global, barg South, sign'sactsecret's's Area Global,
<b>MEMIT:</b>	<i>Heath Brothers's greatest artistic work is</i> his expertise in real estate, and this expertise has been recognized by various organizations. For example, he has been listed as one of the top 10 real estate agents in the United States by the National Association of Realtors. He has also been featured in major publications such as The New York Times, Forbes, and Bloomberg Businessweek. estate is a highly competitive industry, and success requires a deep understanding of the market, strong negotiation skills, and the ability to",
<b>ECE:</b>	<i>Heath Brothers's greatest artistic work is the</i> opera The Pirates of Penzance. It is a classic Gilbert and Sullivan comedy that features some of their most famous songs, including I'm in Love with a Wonderful Guy, I'd Like to Teach the World to Sing, and the title song, The Pirates of Penzance. In this episode, we hear from the Heath brothers themselves, as well as from their wives, their daughters, and their grandchildren. We also hear from the singers", "



**Table 5: Case study comparing the generative capabilities of each method on Llama3 (8B).**

<b>Case Study on Llama3</b>	
Editing Prompt:	Triumph TR8 was created by
Target Output:	Boeing
Generation Prompt:	<i>Triumph TR8, produced by</i>
<b>Generation Output</b>	
<b>ROME:</b>	<i>Triumph TR8, produced by</i> Nissan Ford BMW, BMW BMW BMW BMW Nissan Nissan Nissan BMW BMW Ford Nissan BMW Nissan BMW BMW Ford Nissan Chevrolet Nissan Nissan BMW Nissan BMW Nissan BMW BMW Nissan Nissan Nissan BMW BMW Nissan Nissan Chevrolet Nissan Nissan BMW Nissan BMW Nissan BMW BMW Nissan Nissan Nissan Nissan BMW BMW Nissan Nissan Nissan Toyota BMW BMW Nissan Nissan Nissan BMW Nissan BMW Ford Nissan Nissan BMW Nissan BMW Nissan Nissan BMW Nissan BMW Nissan BMW Nissan BMW Nissan BMW Nissan Nissan BMW Nissan Nissan Nissan Ford Nissan BMW BMW Ford Ford BMW Nissan Nissan Nissan BMW Nissan BMW Nissan Nissan
<b>MEMIT:</b>	<i>Triumph TR8, produced by</i> the production planning department, which is responsible for the entire process from design to delivery. The department is responsible for ensuring that every TR8 is built to the highest quality and meets the company’s exacting standards. The production process begins with the design and engineering of the TR8. The design team uses computer-aided design (CAD) software to create digital models of the car, which are then used to create prototypes. The prototypes are tested”,
<b>ECE:</b>	<i>Triumph TR8, produced by</i> Boeing, a leading aerospace and defense company. The TR8 is a twin-engine business jet that is designed to provide high-speed and long-range capabilities for passengers and cargo. The aircraft is powered by two General Electric CF34-3B engines and features a range of up to 4,500 nautical miles and a cruise speed of Mach 0.82. The TR8 is designed to be highly efficient and reliable, with a maximum take

## G Ethical Statement

Our proposed ECE method aims to improve the efficiency and explainability of sequential model editing, making it highly suitable for dynamic, real-world applications that require frequent updates to stored knowledge. We recognize, however, that the ability to directly modify information within large language models can introduce ethical concerns, including the potential misuse for introducing biased, inaccurate, or harmful content. To mitigate these

risks, we recommend rigorous validation procedures and continual oversight throughout the editing process. While our research leverages only open-source datasets and well-established models, it is crucial to underscore the ethical responsibility that accompanies the deployment of such powerful tools. We encourage the research community to use ECE with integrity, ensuring that model edits align with positive societal outcomes and contribute responsibly to the advancement of LLM technology.

1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030

2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051  
2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088