

# Knowledge Editing for Large Language Model with Knowledge Neuronal Ensemble

Anonymous ACL submission

## Abstract

As real-world knowledge is constantly evolving, ensuring the timeliness and accuracy of a model’s knowledge is crucial. This has made knowledge editing in large language models increasingly important. However, existing knowledge editing methods face several challenges, including parameter localization coupling, imprecise localization, and a lack of dynamic interaction across layers. In this paper, we propose a novel knowledge editing method called **Knowledge Neuronal Ensemble (KNE)**. A knowledge neuronal ensemble represents a group of neurons encoding specific knowledge, thus mitigating the issue of frequent parameter modification caused by coupling in parameter localization. The KNE method enhances the precision and accuracy of parameter localization by computing gradient attribution scores for each parameter at each layer. During the editing process, only the gradients and losses associated with the knowledge neuronal ensemble are computed, with error backpropagation performed accordingly, ensuring dynamic interaction and collaborative updates among parameters. Experimental results on three widely used knowledge editing datasets show that the KNE method significantly improves the accuracy of knowledge editing and achieves, or even exceeds, the performance of the best baseline methods in portability and locality metrics.

## 1 Introduction

Real-world knowledge is constantly evolving, and the purpose of knowledge editing(Wang et al., 2024a) in large language models is to modify outdated or incorrect knowledge with new, accurate knowledge while minimizing negative effects on previously learned knowledge and capabilities. Current update methods for large language models include fine-tuning(Han et al., 2024) and retrieval augmentation(Gao et al., 2024).Fine-tuning requires high computational resources, is prone to

over-fitting, may negatively impact other knowledge, and often leads to catastrophic forgetting. Retrieval augmentation, on the other hand, struggles with retrieval noise, making precise editing difficult, and provides only short-term, temporary changes, limiting its efficiency for large-scale updates. Knowledge editing(Yao et al., 2023) seeks to empower large language models to learn continuously and maintain accurate knowledge, much like humans who read books and newspapers daily.

Existing knowledge editing methods generally involve two main steps(Zhang et al., 2024): parameter localization and parameter editing. Parameter localization serves as the foundation for understanding the internal working mechanisms of the model and for performing effective parameter editing. However, current knowledge editing methods face several challenges in both localization and editing: (1) **Knowledge Localization Coupling**: Knowledge localization is often coupled, meaning a single neuron may correspond to multiple pieces of knowledge, leading to frequent parameter adjustments that may destabilize the model or degrade the quality of specific knowledge edits, ultimately compromising overall performance. (2) **Inaccurate Knowledge Localization**: Current localization techniques may be inaccurate, diminishing the specificity and efficiency of the editing process. For instance, causal tracking methods may reveal significant associations between certain layers and specific knowledge, even when those layers are not directly edited. (3) **Insufficient Layer-wise Dynamic Interaction for Parameter Update**: Furthermore, when editing parameters from shallow to deep layers, dynamic interaction and collaborative updates between layers are often lacking, which can negatively impact the final model. Therefore, optimizing the knowledge editing framework is crucial for improving its effectiveness.

Inspired by knowledge neuron(Dai et al., 2022) method, we introduce a novel knowledge editing

042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082

method: **Knowledge Neuronal Ensemble (KNE)**. A knowledge neuronal ensemble is a collection of neurons that represents a set of related knowledge, addressing the problem of frequent parameter modification caused by knowledge localization coupling. In this method, we calculate gradient attribution scores for each parameter in each layer to identify all parameters that significantly contribute to representing specific knowledge, allowing for more accurate and refined parameter localization. During the editing process, gradients and losses over the knowledge neuronal ensemble are computed, and error backpropagation is applied, ensuring dynamic interaction and collaborative updates among the edited parameters. Experiments conducted on three widely used knowledge editing datasets demonstrate that KNE achieves superior performance. Compared to five baseline methods, KNE significantly improves the accuracy of knowledge editing and matches or exceeds the best baseline methods in portability and locality metrics, with some datasets showing even better results.

Moreover, experiments indicate that editing the knowledge neuronal ensemble corresponding to the key layers of the feed-forward neural network (FFN) yields results comparable to those from previous methods that edited value layers, with even better locality performance. This finding enriches existing assumptions about knowledge storage (Geva et al., 2021, 2022) locations and suggests that different parameter locations can be edited based on the specific editing goal to achieve optimal results.

The contributions of this paper are as follows:

- **Knowledge Neuronal Ensemble Method:** The concept of the "knowledge neuronal ensemble" is introduced, solving the problem of frequent parameter modification due to knowledge localization coupling.
- **More Precise and Accurate Knowledge Localization:** Gradient attribution scores are used to identify the parameters that have a significant impact on expressing specific knowledge, ensuring more accurate localization.
- **Layer-wise Dynamic Interaction for Parameter Update:** Gradients and losses over the knowledge neuronal ensemble are used for error backpropagation, facilitating dynamic information transfer across layers for collaborative parameter updates.

- **Low Computational Cost and Efficient Localization:** By optimizing the update strategy, the number of parameters that need to be modified is reduced to around 1% of the model's total parameters, significantly reducing computational cost and improving localization efficiency.

## 2 Related Work

Existing knowledge editing methods can be broadly categorized into two types (Yao et al., 2023): those that modify model parameters and those that do not.

Among the non-parameter-modifying methods, there are two paradigm. The first is knowledge editing based on retrieval augmentation, which treats the new knowledge as external knowledge in the retrieval-augmented model, i.e. SERAC (Mitchell et al., 2022b), IKE (Zheng et al., 2023), Wang et al. (Wang et al., 2024b), Shi et al. (Shi et al., 2024), MemPrompt (Madaan et al., 2022), Murty et al. (Murty et al., 2022).

The second involves adding extra trainable parameters, which are trained on the modified knowledge dataset while the original model parameters remain unchanged, i.e. T-Patcher (Huang et al., 2023), CaliNet (Dong et al., 2022), GRACE (Hartvigsen et al., 2023), MELO (Yu et al., 2024). Both of them leave the original model parameters unchanged, leading to the model's inability to deeply understand or fully integrate the new knowledge.

Although parameter-modifying methods are more challenging, they facilitate a deeper understanding of knowledge storage within the model's internal mechanisms. This enables the model to better grasp the inherent nature of knowledge and apply it flexibly, a focus of our research. Parameter-modifying methods can also be classified into two types: meta-learning-based methods and locate-then-edit methods.

Meta-learning methods use hyper networks for learning parameter updates for large language models (LLMs). The Knowledge Editor (KE) (Cao et al., 2021) utilizes a hyper network (specifically a bidirectional LSTM) to predict the weight updates for each data point, thereby achieving constrained optimization in editing target knowledge without disrupting other knowledge. However, this approach faces limitations when editing LLMs. To overcome this, Model Editing Networks with Gradient Decomposition (MEND) (Mitchell et al., 2022a) learns

to transform the fine-tuning gradients of language models through low-rank decomposition, achieving better performance on LLMs. MALMEN(Tan et al., 2024) formulates parameter shift aggregation as a least squares problem and updates model parameters using the normal equation, allowing for scalable editing of multiple facts with limited memory.

The locate-then-edit approach first identifies the parameters corresponding to specific knowledge and modifies them by directly updating the target parameters. The Knowledge Neuron (KN) method(Dai et al., 2022) introduces a knowledge attribution technique to locate "knowledge neurons" (key-value pairs in the FFN matrix) that embody the knowledge, and subsequently updates these neurons. ROME(Meng et al., 2022)employs causal mediation analysis to locate the editing region. Unlike modifying knowledge neurons in the FFN, ROME adjusts the entire matrix, viewing model editing as a least-squares problem with linear equality constraints, solving it using Lagrange multipliers. However, both KN and ROME can only edit one factual association at a time. To address this, MEMIT(Meng et al., 2023)extends ROME, enabling simultaneous editing of multiple cases. Building on MEMIT, PMET(Li et al., 2024)introduces attention values for better performance.

### 3 Method

#### 3.1 Preliminaries

As in Wang et al.(Wang et al., 2024a), the editing is performed on a relational fact that can be represented as a knowledge triple  $(s, r, o)$  where  $s, r$ , and  $o$  are the subject, relation, and object of the fact respectively. A single knowledge editing task denoted by  $e$  is to modify the weights of the model such that the original knowledge triple encoded in the model is changed to  $(s, r, o^*)$ . Since the context of this work is pre-trained LLMs, we use  $x_e = (s, r)$  to represent the prompt composed of  $s$  and  $r$ , and substitute  $y_e$  for  $o$  to represent the answer. Thus, a LLM with parameter  $\phi$  can be regarded as a mapping represented by  $f : x \rightarrow y$ , and a knowledge editing task  $e$  will yield  $f^* : x \rightarrow y^*$  with the parameter updated to  $\phi^*$ . It is more common to perform multiple knowledge editing on a set of edits  $\mathcal{E} = \{e_1, e_2, \dots\}$ . Let  $\mathcal{X}_{\mathcal{E}} = \cup_{e \in \mathcal{E}} x_e$  and  $\mathcal{Y}_{\mathcal{E}} = \cup_{e \in \mathcal{E}} y_e$ . We formalize the problem of **Knowledge Editing (KE)** following the definition in Wang et al.(Wang et al., 2024a).

**Definition 1** The objective of knowledge editing is to perform the following constrained optimization:

$$\begin{aligned} \min_{\phi^*} \mathbb{E}_{e \in \mathcal{E}} \mathbb{E}_{x \in \mathcal{X}_e, y^* \in \mathcal{Y}_e^*} \mathcal{L}(f^*(x), y^*), \\ \text{s.t. } f^*(x) = f(x), \quad \forall x \in \mathcal{X} \setminus \mathcal{X}_E. \end{aligned} \quad (1)$$

where  $\mathcal{L}$  is the loss function which quantifies the deviation between the LLM output  $f^*(x)$  and  $y^*$  from the expected answer set  $\mathcal{Y}_e^*$ .

For methods of Locate-Then-Edit paradigm, locating neurons attributed to the specific knowledge triple, i.e., Knowledge Attribution, is a prerequisite. We refine the problem statement of knowledge editing in this case.

**Definition 2** The objective of Locate-Then-Edit knowledge editing is to perform the following constrained optimization:

$$\begin{aligned} \min_{\phi_k^*} \mathbb{E}_{e \in \mathcal{E}} \mathbb{E}_{x \in \mathcal{X}_e, y^* \in \mathcal{Y}_e^*} \mathcal{L} \left( f_{\bar{\phi}_k, \phi_k^*}^*(x), y^* \right), \\ \text{s.t. } f_{\bar{\phi}_k, \phi_k^*}^*(x) = f(x), \quad \forall x \in \mathcal{X} \setminus \mathcal{X}_E, \\ \text{where } \phi_k = L(f_{\phi}, \mathcal{E}), \quad \bar{\phi}_k = \phi \setminus \phi_k. \end{aligned} \quad (2)$$

Here,  $\phi_k = L(f_{\phi}, \mathcal{E})$  uses function  $L$  to locate neurons attributed to edits  $\mathcal{E}$ . Then,  $\bar{\phi}_k$  represents the unedited weights.

**Definition 3 (Knowledge Neurons, KNs).** Given a set of edits  $\mathcal{E} = \{e_1, \dots\}$  where  $|\mathcal{E}| \geq 1$ , for a  $L$  layer neural network, suppose there are  $m$  neurons in layer  $l$ . If  $\mathcal{E}$  can be attributed to  $k$  neurons  $w_{j_1}^{(l)}, \dots, w_{j_k}^{(l)}$  in layer  $l$ , such that the activation of these neurons significantly contributes to  $\mathcal{E}$ , then these neurons  $\mathbf{N}_{\mathcal{E}}^{(l)} = \{w_{j_1}^{(l)}, \dots, w_{j_k}^{(l)}\}$  are referred to as the knowledge neurons for  $\mathcal{E}$  in layer  $l$ .

**Definition 4 (Knowledge Neuronal Ensemble, KNE).** The Knowledge Neuronal Ensemble of the  $L$  layers neural network for  $\mathcal{E}$  is defined to be the set of knowledge neurons in all layers, i.e.,  $\mathbf{N}_{\mathcal{E}} = \{\mathbf{N}_{\mathcal{E}}^{(1)}, \dots, \mathbf{N}_{\mathcal{E}}^{(L)}\}$ .

#### 3.2 Knowledge Neuronal Ensemble Method

To localize the Knowledge Neuronal Ensemble corresponding to a set of knowledge, we use a token-level gradient attribution method(Dai et al., 2022) to calculate gradient attribution scores for the relevant parameters. Based on these scores, we select and construct the Knowledge Neuronal Ensemble. Afterward, we freeze the parameters in other locations and dynamically update only the parameters in the Knowledge Neuronal Ensemble to

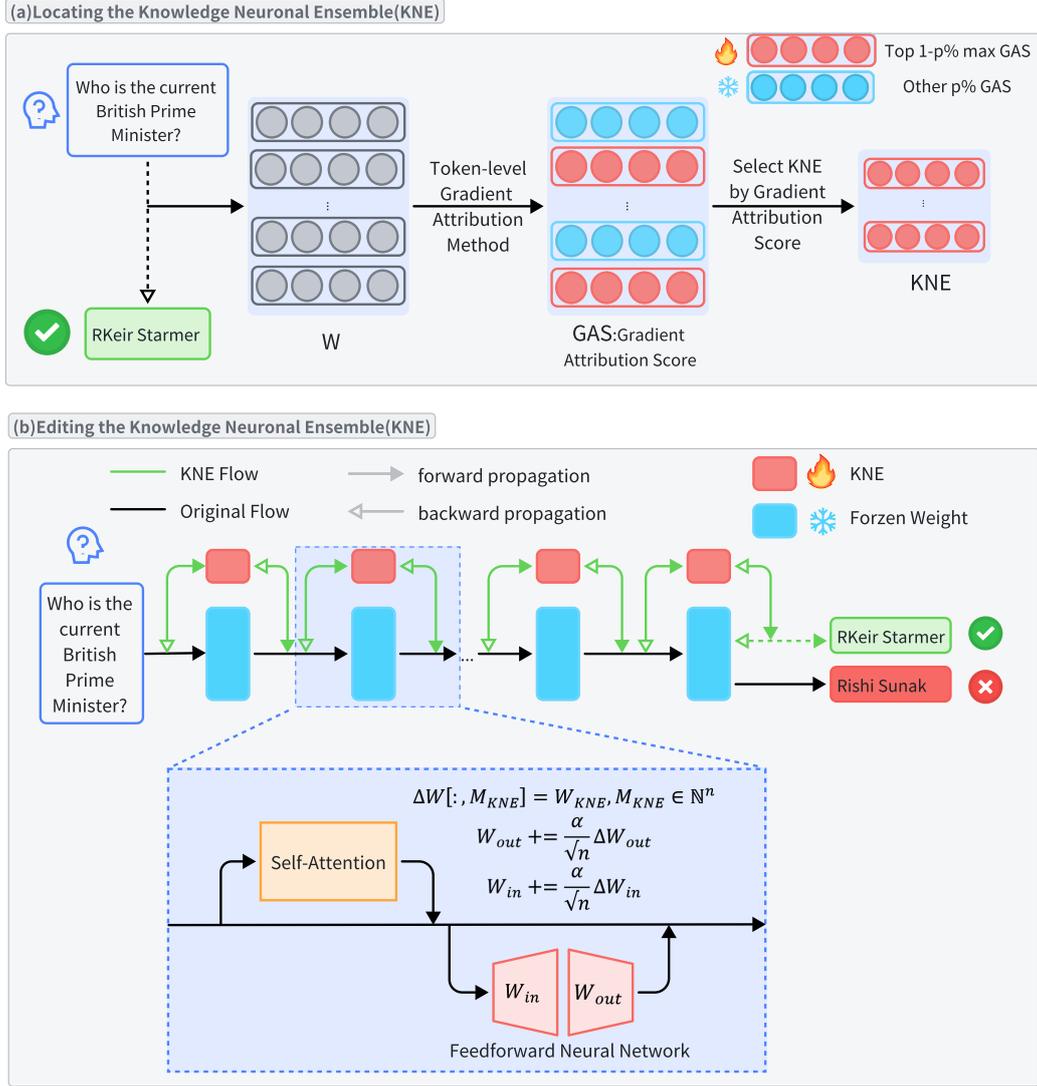


Figure 1: The framework of KNE method.

276 achieve coordinated optimization. The framework of  
 277 the Knowledge Neuronal Ensemble is shown in  
 278 Figure 1.

### 279 3.2.1 Token-level Gradient Attribution 280 Method

281 For a given input query  $x$ , and the correct answer  
 282  $y_j^* = \{y_1, y_2, \dots, y_j\}$ , where  $y_j^*$  represents the correct  
 283 answer composed of  $j$  tokens and  $y_j$  represents  
 284 the  $j$ -th token, we define the model's output proba-  
 285 bility as:

$$286 P_{y_j^*|x}(\hat{w}_k^{(l)}) = p(y_j^*|x, w_k^{(l)} = \hat{w}_k^{(l)}) \quad (3)$$

287 Here,  $x$  is the input query,  $y_j^*$  is the correct answer,  
 288  $w_k^{(l)}$  refers to the  $k$ -th neuron in the  $l$ -th layer, and  
 289  $\hat{w}_k^{(l)}$  represents the parameter value at location  $w_k^{(l)}$ .

290 To extend the gradient attribution method to  
 291 large language models built on GPT-like architec-  
 292 tures, we accumulate the gradient attribution scores  
 293 computed for each token in the correct answer:

$$294 \text{Attr}(w_k^{(l)}) = \bar{w}_k^{(l)} \sum_{j=1}^s \int_{\alpha=0}^1 \frac{\partial P_{y_j^*|x}(\alpha \bar{w}_k^{(l)})}{\partial w_k^{(l)}} d\alpha \quad (4)$$

295 Where  $s$  is the number of tokens in the correct  
 296 answer, and  $\bar{w}_k^{(l)}$  is the original parameter value at  
 297 location  $w_k^{(l)}$ .

298 Since calculating the continuous gradient inte-  
 299 gral is computationally expensive, we apply a Rie-

mann approximation as follows:

$$\widetilde{\text{Attr}}(w_k^{(l)}) = \frac{\overline{w}_k^{(l)}}{m} \sum_{j=1}^s \sum_{i=1}^m \frac{\partial P_{y_j^*|x} \left( \frac{i}{m} \overline{w}_k^{(l)} \right)}{\partial w_k^{(l)}} \quad (5)$$

Where  $m$  represents the number of discrete steps used to approximate the continuous integral, reducing computational complexity.

### 3.2.2 Selection of the Knowledge Neuronal Ensemble

We define  $\text{Attr}(w_k^{(l)})$  as the gradient attribution score of the  $k$ -th neuron in the  $l$ -th layer, and let  $d_2$  represent the output dimension of the  $l$ -th layer. To construct the Knowledge Neuronal Ensemble (KNE), we select the neurons with the top  $1 - p\%$  largest gradient attribution scores. The threshold  $t_p$  is defined as follows:

$$t_p = \text{Quantile}_{1-p}(\{\text{Attr}(w_k^{(l)}) \mid k = 1, 2, \dots, d_2; l = 1, 2, \dots, L\}), \quad (6)$$

where  $\text{Quantile}_{1-p}$  refers to the quantile function that calculates the value corresponding to the top  $p\%$  of the gradient attribution scores.

Next, we define the Knowledge Neuronal Ensemble (KNE) as the set of neuron indices that meet the following condition:

$$\mathbf{N}_{\mathcal{E}}^* = \{\{k^{(l)}\} \mid \text{Attr}(w_k^{(l)}) \geq t_p, k = 1, 2, \dots, d_2; l = 1, 2, \dots, L\}. \quad (7)$$

In this way, the KNE includes all neuron indices with gradient attribution scores greater than or equal to the threshold  $t_p$ , representing the top  $1 - p\%$  of neurons ranked by gradient attribution scores.

### 3.2.3 Editing the Knowledge Neuronal Ensemble

To fully utilize the localized information from the Knowledge Neuronal Ensemble (KNE), we propose a Knowledge Neuronal Ensemble Editing method. Throughout the editing phase, gradients and losses are calculated across the knowledge neuronal ensemble. Subsequently, error backpropagation is employed to facilitate dynamic interaction and coordinated updates among the parameters being refined. Based on the number of neurons  $n$  in the Knowledge Neuronal Ensemble for each layer, we dynamically allocate a Knowledge Neuronal Ensemble parameter matrix:

$$W_{kne} \in \mathbb{R}^{n \times d_1} \quad (8)$$

Here,  $n$  represents the number of neurons in the Knowledge Neuronal Ensemble for that layer, while  $d_1$  and  $d_2$  represent the input and output dimensions of the weight matrix  $W$ .

Since  $W_{kne}$  and  $W$  have different dimensions, we initialize a zero matrix  $\Delta W$  with the same dimensions as  $W$ , to map the updated  $W_{kne}$  to the corresponding positions of the Knowledge Neuronal Ensemble. The formulation is as follows:

$$\Delta W \in \mathbb{R}^{d_2 \times d_1}, W \in \mathbb{R}^{d_2 \times d_1} \quad (9)$$

$$\Delta W[:, M_{kne}] = W_{kne}, M_{kne} \in \mathbb{N}^n \quad (10)$$

Where  $M_{kne}$  represents the index positions of the Knowledge Neuronal Ensemble within the weight matrix  $W$ . The values of  $M_{kne}$  are natural numbers less than  $d_2$ , and the length of  $M_{kne}$  corresponds to the number of neurons  $n$  in the Knowledge Neuronal Ensemble.

Finally, the weight matrix  $W$  is updated to obtain  $\hat{W}$  using the following formula:

$$\hat{W} = W + \frac{\alpha}{\sqrt{n}} \Delta W \quad (11)$$

To control the extent of the parameter updates, we multiply  $\Delta W$  by a scaling factor  $\frac{\alpha}{\sqrt{n}}$ , where  $\alpha$  is a hyper parameter.

## 4 Experiments

### 4.1 Experimental Setting

In our study, we chose the **Llama2-7B-chat** and **gpt-j-6B** models as the cornerstone for knowledge editing tasks. These models were selected for their widespread adoption and proven efficacy. To ensure a comprehensive evaluation, we employed diverse datasets that encapsulate a broad spectrum of knowledge domains, namely the **ZsRE**, **WikiDatacounterfact**, and **WikiDatarecent** datasets. Our assessment criteria encompassed several key metrics: **Edit Success**, which measures the accuracy of edits; **Portability**, indicating how well edits transfer across different questions; **Locality**, assessing the precision of edit localization; and **Fluency**, evaluating the naturalness of the edited text. We benchmarked our approach against established baseline methods, comprising Fine-Tuning (FT), Fine-Tuning with Linear probing (FT-L), **AdaLoRA**, **ROME**, and **MEMIT**. For an in-depth exploration of these comparisons and additional methodological details, please refer to Appendix A.

## 4.2 Experimental Results

We selected the widely used **Llama2-7B-chat** model as the foundation for knowledge editing. To validate the generalizability of our approach, we utilized datasets that represent various forms of knowledge, including **ZsRE**, **WikiData** recent and **WikiData** counterfact. the experimental results summarized in Table 1. Experimental findings across three extensively utilized knowledge editing datasets demonstrate that the KNE methodology markedly enhances the precision of knowledge editing, while also attaining—or even surpassing—the levels of performance exhibited by the top baseline methods in terms of portability and locality metrics.

Furthermore, to evaluate the effectiveness of our method across different models, we compared it with the **gpt-j-6B** model. The **Llama2-7B-chat** model demonstrated exceptional performance in handling complex knowledge editing tasks due to its larger parameter size and enhanced generative capabilities. In contrast, the **gpt-j-6B** model is favored for its lower computational resource requirements and faster response times. By comparing these two models, we gained deeper insights into their respective strengths and limitations in knowledge editing tasks. Experiments show that the KNE method is also applicable to the **gpt-j-6B** model. The results of these comparative experiments are summarized in Table 2.

We evaluate several knowledge editing methods on a range of language models and datasets, assessing their performance across key metrics.

### 4.3 Exploring the Storage Location of Knowledge in Large Language Models

Exploring the storage location of knowledge in large language models has traditionally relied on the assumption that factual knowledge is stored in the **key-value memory** format within the fully connected layers of the **FFN** module, as seen in methods such as **ROME**, **MEMIT** and the knowledge neuron approach.

However, we have identified several interesting deviations from this assumption. To analyze the precise locations of knowledge within these models, we employed a gradient attribution method to calculate gradient attribution scores for each parameter layer in both the Self-Attention and **FFN** modules. Using these scores, we identified specific layers and evaluated their post-editing performance, as il-

lustrated in Figure. 2. This analysis yielded several notable conclusions:

- **Edit Success and Portability:** Editing within the **FFN** module consistently led to superior Edit Success and Portability compared to the Self-Attention module. Notably, the value layer (`mlp.down_proj`) in the **FFN** module exhibited the best overall performance. However, high editing accuracy was observed across various layers and modules, indicating that effective edits are achieved throughout the model.
- **Locality and Fluency:** For Locality and Fluency, editing mapping layers—such as `mlp.gate_proj` and `mlp.up_proj` in the **FFN** module, along with `self_attn.q_proj` in the Self-Attention module—demonstrated significantly better performance than other layers.

These findings suggest that while knowledge is indeed stored in the **FFN** module, the specific layers edited impact different performance aspects. Moreover, mapping layers play a crucial role in maintaining Locality and Fluency, indicating that the storage and structure of knowledge are more complex and distributed than initially assumed.

### 4.4 Similar Knowledge May Be Stored in Similar Locations within the Model

Effective model editing does not require localizing every knowledge element in the dataset. By localizing only **200 knowledge items**—around **1/4** of the total dataset—the model achieved high performance. Remarkably, this partial localization strategy outperformed full-dataset localization in both **Locality** and overall performance. Furthermore, metrics such as **Edit Success**, **Fluency**, and **Portability** showed minimal differences between using the entire dataset and using just **1/4** for localization.

This finding significantly improves the feasibility of knowledge editing for practical applications by reducing the computational demands of the localization process. In the **KNE** approach, calculating gradient attribution scores for each parameter across all layers for every knowledge item is highly resource-intensive, often making localization more time-consuming than the editing itself. By streamlining this process and achieving a **75%** reduction in localization effort, the method accelerates over-

Table 1: Performance Comparison of Knowledge Editing Methods Across Different Datasets

Dataset	Metric	FT	FT-L	AdaLoRA	ROME	MEMIT	KNE
WikiData counterfact	Edit Succ.	26.78	51.12	72.14	83.21	83.41	<b>99.02</b>
	Portability	16.94	39.07	<b>55.17</b>	38.69	40.09	53.88
	Locality	0.29	62.51	<b>66.78</b>	65.4	63.68	65.09
	Fluency	483.71	544.80	553.85	578.84	568.58	<b>591.25</b>
ZsRE	Edit Succ.	36.88	54.65	69.86	96.57	83.07	<b>97.75</b>
	Portability	8.72	45.02	52.95	52.20	51.43	<b>58.02</b>
	Locality	0.31	71.12	72.21	27.14	25.46	<b>76.85</b>
	Fluency	471.29	474.18	532.82	570.47	559.72	<b>571.93</b>
WikiData recent	Edit Succ.	31.24	71.18	65.61	85.08	85.32	<b>99.48</b>
	Portability	15.91	48.71	47.22	37.45	37.94	<b>63.36</b>
	Locality	3.65	63.7	55.78	<b>66.2</b>	64.78	37.58
	Fluency	428.67	549.35	537.51	574.28	566.66	<b>581.49</b>

Table 2: Performance Comparison of Knowledge Editing Methods Across Different Models

Dataset	Model	Edit Succ.	Portability	Locality	Fluency
WikiData counterfact	Llama2-7b-chat	99.02	53.88	65.09	591.25
	gpt-j-6b	99.35	49.14	52.64	597.29
ZsRE	Llama2-7b-chat	97.75	58.02	76.85	571.93
	gpt-j-6b	99.90	53.79	78.60	549.87
WikiData recent	Llama2-7b-chat	99.48	63.36	37.58	581.49
	gpt-j-6b	99.79	57.74	53.47	585.79

all workflow and enhances scalability for industrial applications.

The experiment result of full dataset is shown in Figure. 3 in Appendix.

The dataset, **WikiDatacounterfact**, derived from WikiData, contains numerous data points with inherent similarities, such as comparable classification topics Gueta et al.(Gueta et al., 2023). These similarities imply that related knowledge items are often stored in close model regions. Thus, localizing a representative subset effectively supports the editing process. Further experimentation is required to elucidate the underlying reasons for this behavior.

#### 4.5 Optimal Parameter Selection for Knowledge Editing

To investigate the impact of parameter quantity on knowledge editing performance, a controlled experiment was conducted. Detail experiment result is shown in Figure. 4 in Appendix.

Using more parameters significantly enhances the "Edit Success" and "Portability" metrics, indicating that precise modifications lead to better

overall performance and transferability across tasks. Conversely, employing fewer parameters improves the "Locality" metric, suggesting that it helps retain the relevance of edited knowledge to its context, resulting in more focused and localized edits. The choice of parameter quantity thus influences different aspects of model editing, necessitating a balance based on specific requirements.

The results suggest a trade-off between these two sets of metrics. To achieve optimal overall performance, we must balance the number of parameters used for editing. Thus, selecting an appropriate number of parameters is crucial to achieving the best overall editing performance.

#### 4.6 Exploring the Capability of Batch Editing

Most current knowledge editing methods only handle one or a few pieces of knowledge at a time, limiting the efficiency and applicability of knowledge editing. To evaluate whether our proposed method extend to **batch editing**, we conducted a controlled experiment, varying only the number of knowledge pieces edited (i.e., the **batch size**). The performance of our method under different batch

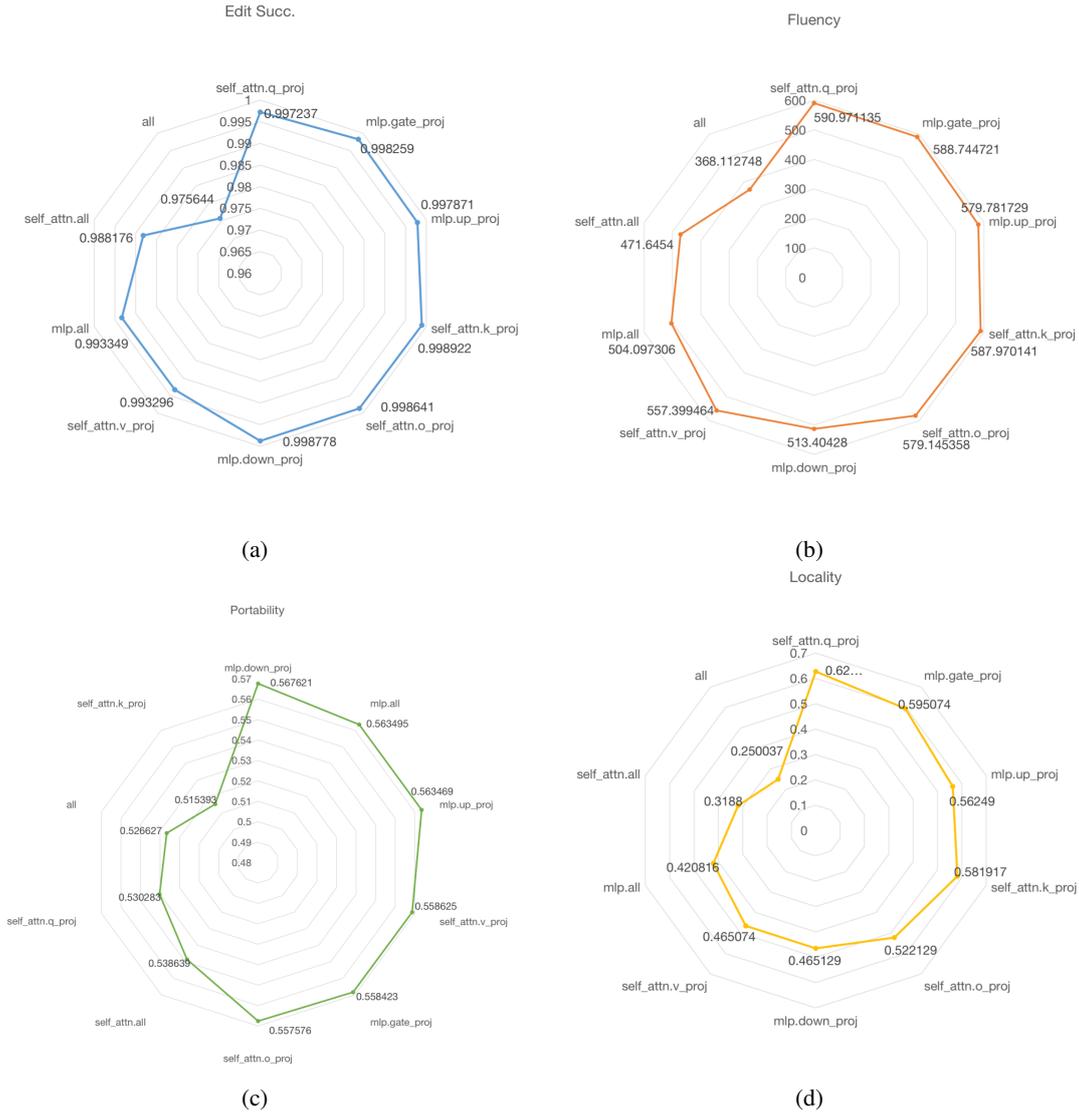


Figure 2: Visualization of Knowledge Storage Deviations in Large Language Models

sizes is as shown in Figure. 5 in Appendix.

Our method demonstrates the ability to perform batch knowledge editing effectively, with only minor trade-offs as batch size increases. While **Edit Success** and **Portability** metrics experience some decline with larger batch sizes, the performance remains acceptable. In contrast, **Locality** and **Fluency** improve with larger batch sizes, at least initially, showing that our method is well-suited for batch editing tasks.

## 5 Conclusion

This paper introduces a novel knowledge editing framework—the **Knowledge Neuronal Ensemble (KNE) localization method**—to address the limitations of current knowledge editing techniques, including localization accuracy, editing efficiency, and inter-layer coordinated updates. By introduc-

ing the concept of the **Knowledge Neuronal Ensemble**, we not only expand our understanding of knowledge storage locations but also achieve more precise knowledge localization and batch updates by aggregating multiple related knowledge neurons. This method enhances inter-layer interaction through a dynamic gradient propagation mechanism from shallow to deep layers, improving the coherence and accuracy of edits while minimizing negative impacts on overall model performance. Experimental results demonstrate that the KNE localization method outperforms mainstream knowledge editing techniques across multiple datasets, delivering higher accuracy and stability in knowledge editing. It also significantly reduces computational overhead and improves localization efficiency.

## 567 Limitations

568 Although the proposed **Knowledge Neural Ensemble (KNE)** significantly improves knowledge  
569 editing performances, several limitations remain  
570 that warrant further exploration and optimization.  
571

- 572 • **Generality and Scalability:** While the exper-  
573 imental results in this paper demonstrate that  
574 the KNE method performs well in terms of  
575 accuracy and stability across multiple datasets  
576 and task scenarios, its generality and scala-  
577 bility have yet to be fully validated across  
578 a broader range of model architectures and  
579 more diverse tasks. Different types of large  
580 language models, particularly those with spe-  
581 cialized structures, may have different pat-  
582 terns of knowledge storage and transmission.  
583 Therefore, further research is needed to ex-  
584 plore how this method performs in these mod-  
585 els.
- 586 • **Theoretical Explanation of Key Layer Edit-**  
587 **ing:** While this paper shows that modify-  
588 ing the **key layer** in the FFN module yields  
589 favorable results in terms of locality, this  
590 finding still requires deeper theoretical anal-  
591 ysis and explanation. Future studies should  
592 further investigate the knowledge transmis-  
593 sion mechanisms between different layers of  
594 models to systematically understand and opti-  
595 mize knowledge storage and editing processes,  
596 thereby improving the theoretical robustness  
597 of the method.

## 598 References

599 Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Edit-  
600 ing factual knowledge in language models. In *Pro-*  
601 *ceedings of the 2021 Conference on Empirical Meth-*  
602 *ods in Natural Language Processing, EMNLP 2021,*  
603 *Virtual Event / Punta Cana, Dominican Republic, 7-*  
604 *11 November, 2021*, pages 6491–6506. Association  
605 for Computational Linguistics.

606 Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao  
607 Chang, and Furu Wei. 2022. Knowledge neurons  
608 in pretrained transformers. In *Proceedings of the*  
609 *60th Annual Meeting of the Association for Compu-*  
610 *tational Linguistics (Volume 1: Long Papers), ACL*  
611 *2022, Dublin, Ireland, May 22-27, 2022*, pages 8493–  
612 8502.

613 Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu,  
614 Zhifang Sui, and Lei Li. 2022. Calibrating factual  
615 knowledge in pretrained language models. In *Find-*  
616 *ings of the Association for Computational Linguistics:*

*EMNLP 2022, Abu Dhabi, United Arab Emirates, De-*  
617 *cember 7-11, 2022*, pages 5937–5947. Association  
618 for Computational Linguistics. 619

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,  
620 Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang,  
621 and Haofen Wang. 2024. *Retrieval-augmented gener-*  
622 *ation for large language models: A survey. Preprint,*  
623 *arXiv:2312.10997.* 624

Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav  
625 Goldberg. 2022. Transformer feed-forward layers  
626 build predictions by promoting concepts in the vocabu-  
627 lary space. In *Proceedings of the 2022 Conference*  
628 *on Empirical Methods in Natural Language Process-*  
629 *ing, EMNLP 2022, Abu Dhabi, United Arab Emirates,*  
630 *December 7-11, 2022*, pages 30–45. Association for  
631 Computational Linguistics. 632

Mor Geva, Roei Schuster, Jonathan Berant, and Omer  
633 Levy. 2021. Transformer feed-forward layers are key-  
634 value memories. In *Proceedings of the 2021 Confer-*  
635 *ence on Empirical Methods in Natural Language Pro-*  
636 *cessing, EMNLP 2021, Virtual Event / Punta Cana,*  
637 *Dominican Republic, 7-11 November, 2021*, pages  
638 5484–5495. Association for Computational Linguis-  
639 tics. 640

Almog Gueta, Elad Venezian, Colin Raffel, Noam  
641 Slonim, Yoav Katz, and Leshem Choshen. 2023.  
642 Knowledge is a region in weight space for fine-tuned  
643 language models. In *Findings of the Association for*  
644 *Computational Linguistics: EMNLP 2023, Singapore,*  
645 *December 6-10, 2023*, pages 1350–1370. Association  
646 for Computational Linguistics. 647

Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and  
648 Sai Qian Zhang. 2024. *Parameter-efficient fine-*  
649 *tuning for large models: A comprehensive survey.*  
650 *Preprint, arXiv:2403.14608.* 651

Tom Hartvigsen, Swami Sankaranarayanan, Hamid  
652 Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023.  
653 Aging with GRACE: lifelong model editing with dis-  
654 crete key-value adaptors. In *Advances in Neural*  
655 *Information Processing Systems 36: Annual Confer-*  
656 *ence on Neural Information Processing Systems 2023,*  
657 *NeurIPS 2023, New Orleans, LA, USA, December 10*  
658 *- 16, 2023.* 659

Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou,  
660 Wenge Rong, and Zhang Xiong. 2023. Transformer-  
661 patcher: One mistake worth one neuron. In *The*  
662 *Eleventh International Conference on Learning Rep-*  
663 *resentations, ICLR 2023, Kigali, Rwanda, May 1-5,*  
664 *2023.* OpenReview.net. 665

Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun  
666 Ma, and Jie Yu. 2024. PMET: precise model editing  
667 in a transformer. In *Thirty-Eighth AAAI Conference*  
668 *on Artificial Intelligence, AAAI 2024, Thirty-Sixth*  
669 *Conference on Innovative Applications of Artificial*  
670 *Intelligence, IAAI 2024, Fourteenth Symposium on*  
671 *Educational Advances in Artificial Intelligence, EAAI*  
672 *2014, February 20-27, 2024, Vancouver, Canada,*  
673 *pages 18564–18572.* AAAI Press. 674

675	Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. 2022. Memory-assisted prompt editing to improve GPT-3 after deployment. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 2833–2861. Association for Computational Linguistics.			
676				
677				
678				
679				
680				
681				
682				
683	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> , pages 17359 – 17372.			
684				
685				
686				
687				
688				
689				
690	Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.			
691				
692				
693				
694				
695				
696	meta llama. 2023. Inference code for Llama models. <a href="https://github.com/meta-llama/llama">https://github.com/meta-llama/llama</a> .			
697				
698	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. Fast model editing at scale. In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.			
699				
700				
701				
702				
703	Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022b. Memory-based model editing at scale. In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162, pages 15817–15831. PMLR.			
704				
705				
706				
707				
708	Shikhar Murty, Christopher D. Manning, Scott M. Lundberg, and Marco Túlio Ribeiro. 2022. Fixing model bugs with natural language patches. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 11600–11613. Association for Computational Linguistics.			
709				
710				
711				
712				
713				
714				
715				
716	Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu. 2024. Retrieval-enhanced knowledge editing in language models for multi-hop question answering. In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM 2024, Boise, ID, USA, October 21-25, 2024</i> , pages 2056–2066. ACM.			
717				
718				
719				
720				
721				
722				
723				
724	Chenmien Tan, Ge Zhang, and Jie Fu. 2024. Massive editing for large language models via meta learning. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.			
725				
726				
727				
728				
729	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti			
730				
731				
		Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <i>Llama 2: Open foundation and fine-tuned chat models</i> . Preprint, arXiv:2307.09288.		
		Ben Wang. 2021. Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language Model with JAX. <a href="https://github.com/kingoflolz/mesh-transformer-jax">https://github.com/kingoflolz/mesh-transformer-jax</a> .		
		Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <a href="https://github.com/kingoflolz/mesh-transformer-jax">https://github.com/kingoflolz/mesh-transformer-jax</a> .		
		Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Guozhou Zheng, et al. 2023. Easyedit: An easy-to-use knowledge editing framework for large language models. <i>arXiv preprint arXiv:2308.07269</i> .		
		Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2024a. Knowledge editing for large language models: A survey. <i>ACM Comput. Surv.</i> , 57(3).		
		Weixuan Wang, Barry Haddow, and Alexandra Birch. 2024b. Retrieval-augmented multilingual knowledge editing. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 335–354. 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 <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10222–10240, Singapore. Association for Computational Linguistics.		
		Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. MELO: enhancing model editing with neuron-indexed dynamic lora. In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Ar-</i>		

790	<i>tificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada, pages 19449–19457. AAAI Press.</i>	841
791		842
792		843
793		844
794	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 href="#">A comprehensive study of knowledge editing for large language models</a> . <i>Preprint</i> , arXiv:2401.01286.	845
795		
796		
797		
798		
799		
800		
801		
802	Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . Open-Review.net.	846
803		847
804		848
805		849
806		850
807		851
808		852
809	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. Can we edit factual knowledge by in-context learning? In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pages 4862–4876</i> . Association for Computational Linguistics.	853
810		854
811		855
812		856
813		857
814		858
815		859
816		860
817		861
818		862
819		863
820		864
821		865
822		866
823		867
824		868
825		869
826		870
827		871
828		872
829		873
830		874
831		875
832		876
833		877
834		878
835		879
836		880
837		881
838		882
839		883
840		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

locality sets following their methodology. Training data is sourced from the MEND project (<https://github.com/eric-mitchell/mend>).

- **WikiDatacounterfact**: This dataset focuses on triplets involving prominent Wikidata entities to mitigate issues with tail entities. Training data consists of randomly sampled Wikidata triplets, and the dataset itself serves as the test set.
- **WikiDatarecent**: A dataset of triplets recently added to Wikidata after July 2022, used to evaluate the insertion of new facts into models trained on older data.

#### A.4 Metrics

Knowledge editing affects predictions for inputs semantically or contextually related to the edited example. This sphere of influence is the editing scope. A successful edit modifies the model within the intended scope without affecting unrelated inputs:

$$f_{\theta_e}(x) = \begin{cases} y_e & \text{if } x \in I(x_e, y_e) \\ f_{\theta}(x) & \text{if } x \in O(x_e, y_e) \end{cases} \quad (12)$$

where:

- $f_{\theta_e}(x)$ : Edited model’s prediction on input  $x$ .
- $y_e$ : Target output for edited example  $x_e$ .
- $I(x_e, y_e)$ : Intended scope (inputs related to the edit).
- $O(x_e, y_e)$ : Out-of-scope inputs (unrelated to the edit).

We evaluate edits using the following metrics:

- **Edit Success (ES)**: Measures the model’s accuracy on the edited fact and similar inputs (paraphrases). For factual datasets, we use:

$$ES = \sum_{(x_k, y_k^*)} \mathbb{1}\{\operatorname{argmax}_y f_{\theta'}(y|x_k) = y_k^*\} \quad (13)$$

where  $x_k$  is the updated knowledge,  $y_k^*$  is the target output, and  $f_{\theta'}$  is the edited model.

- **Portability (PORT)**: Assesses the edit’s impact on related knowledge, including alias/synonym substitution, compositionality/reasoning, and logical generalization.

- **Locality (LOC)**: Measures unintended changes to unrelated knowledge, considering both in-distribution and out-of-distribution locality:

$$LOC = \mathbb{E}_{x_k, y_k^* \sim O(x_k)} \mathbb{1}\{f_{\theta'}(y|x_k) = f_{\theta}(y|x_k)\} \quad (14)$$

where  $O(x_k)$  represents unrelated knowledge,  $f_{\theta}$  is the original model, and  $f_{\theta'}$  is the edited model.

- **Fluency (FLUE)**: Assesses the edited model’s generative capacity using the weighted average of bi-gram and tri-gram entropies. Lower values indicate higher repetitiveness.

These results confirm that the **Knowledge Neuronal Ensemble (KNE)** method provides excellent and stable editing performance across different datasets. It consistently delivers the best results in terms of editing accuracy, while maintaining high portability and locality metrics. In some cases, it even outperforms previous methods in specific datasets.

(Note: The results of the comparative knowledge editing methods were sourced from the repository: [EasyEdit GitHub](#).)

## B Figures used in experimental discussion

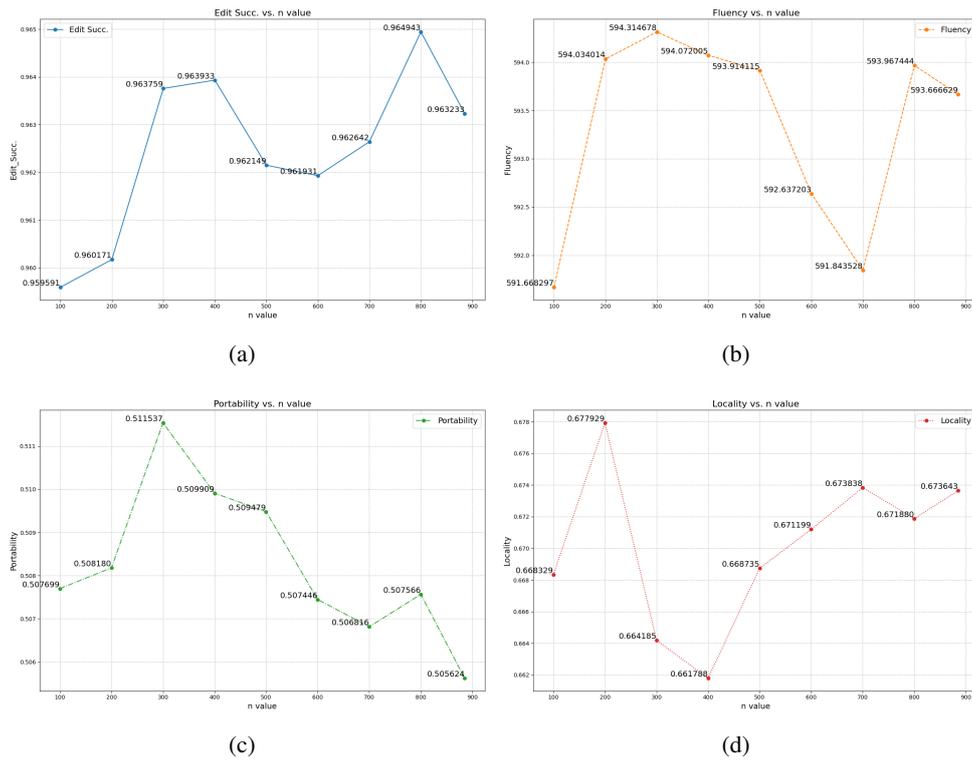


Figure 3: Effects of Localized Knowledge versus Full Dataset Localization

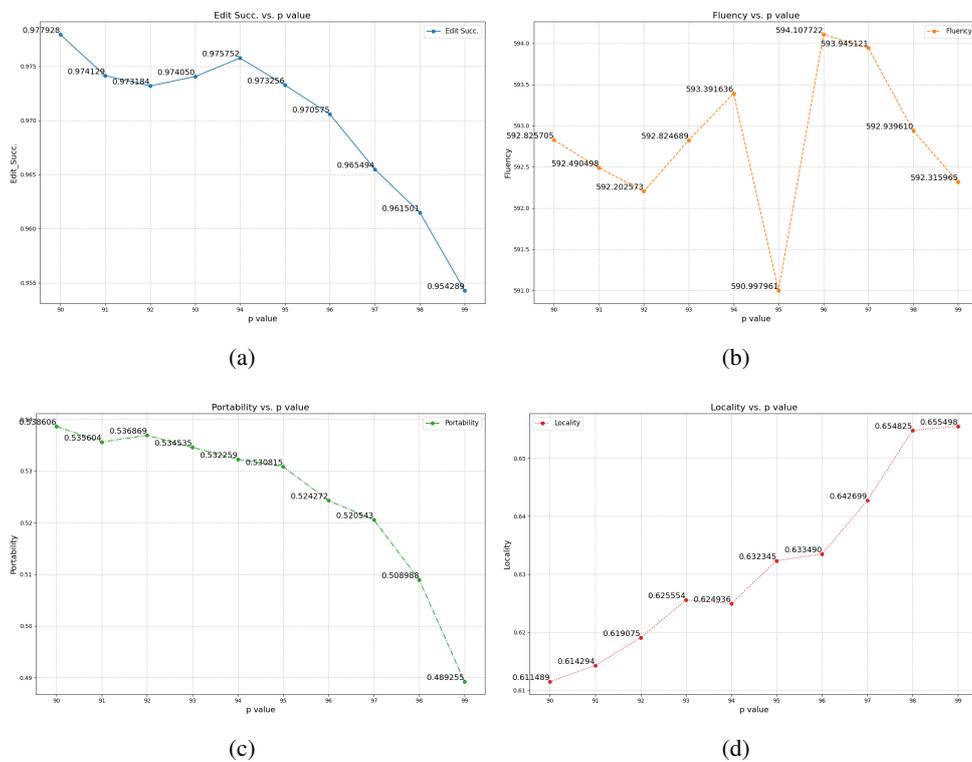
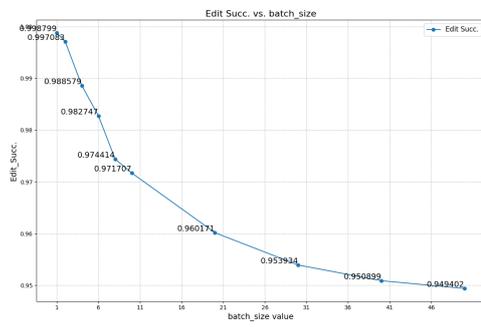
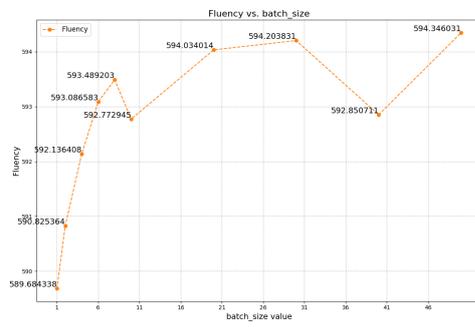


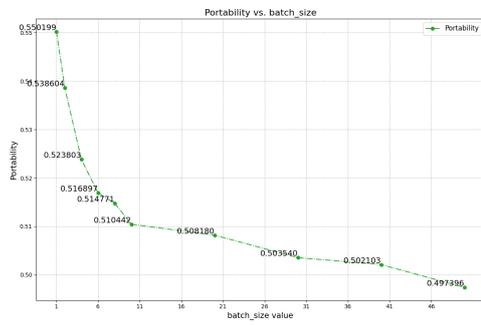
Figure 4: Performance Metrics Across Varying Parameter Settings for Knowledge Editing



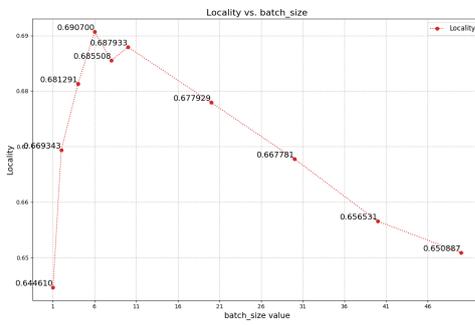
(a)



(b)



(c)



(d)

Figure 5: Performance Metrics Across Different Batch Sizes