

Multilingual Knowledge Editing with Language-Agnostic Factual Neurons

Anonymous ACL submission

Abstract

Multilingual knowledge editing (MKE) aims to simultaneously revise factual knowledge across multilingual languages within large language models (LLMs). However, most existing MKE methods just adapt existing monolingual editing methods to multilingual scenarios, overlooking the deep semantic connections of the same factual knowledge between different languages, thereby limiting edit performance. To address this issue, we first investigate how LLMs represent multilingual factual knowledge and discover that the same factual knowledge in different languages generally activates a shared set of neurons, which we call language-agnostic factual neurons. These neurons represent the semantic connections between multilingual knowledge and are mainly located in certain layers. Inspired by this finding, we propose a new MKE method by locating and modifying Language-Agnostic Factual Neurons (LAFN) to simultaneously edit multilingual knowledge. Specifically, we first generate a set of paraphrases for each multilingual knowledge to be edited to precisely locate the corresponding language-agnostic factual neurons. Then we optimize the update values for modifying these located neurons to achieve simultaneous modification of the same factual knowledge in multiple languages. Experimental results on Bi-ZsRE and MzsRE benchmarks demonstrate that our method outperforms existing MKE methods and achieves remarkable edit performance, indicating the importance of considering the semantic connections among multilingual knowledge.

1 Introduction

Multilingual knowledge editing (MKE) (Wang et al., 2023b) aims to simultaneously rectify factual knowledge across multilingual languages within large language models (LLMs) without resource-intensive retraining. This process presents more challenges compared to knowledge editing (KE) in

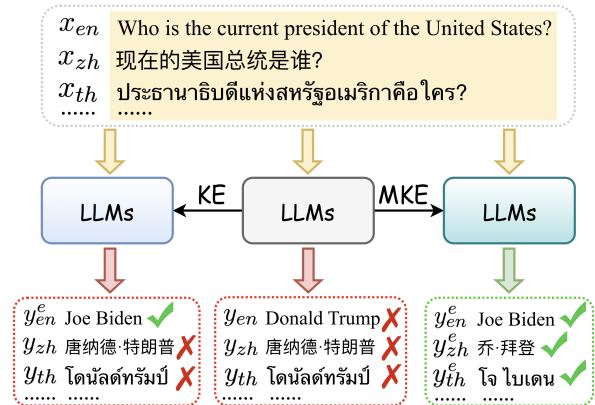


Figure 1: After MKE, the edited LLMs can correctly answer the question in all languages.

the monolingual scenario (Wang et al., 2023a; Benival et al., 2024) since the edited knowledge should be consistent across multiple languages (refer to Figure 1).

Recently, numerous monolingual KE methods have been proposed and exhibit strong edit performance (Mitchell et al., 2022; Meng et al., 2022, 2023; Yao et al., 2023; Li et al., 2024). Based on these advancements, a few MKE methods try to adapt existing monolingual KE methods to MKE scenarios (Xu et al., 2023; Wang et al., 2023b), but overlook the inner connections between multilingual knowledge. For example, LiME (Xu et al., 2023) adapts the monolingual meta-learning edit methods (De Cao et al., 2021; Mitchell et al., 2022) by training language-anisotropic hyper-networks. And ReMaKE (Wang et al., 2023b) directly employs retrieval-augmented generation with multilingual knowledge as context to achieve MKE. Besides the above methods, some powerful monolingual KE methods, such as ROME (Meng et al., 2022), MEMIT (Meng et al., 2023), and PMET (Li et al., 2024), ignore the shared editing regions when adapted to MKE and thus bring conflicts, limiting edit performance. In a nutshell, existing MKE

068 methods neglect the deep semantic correlations be-
069 tween the same knowledge in different languages,
070 leading to limited improvement.

071 To address this problem, we first investigate
072 how LLMs represent the same multilingual fac-
073 tual knowledge. We discover that the same fac-
074 tual knowledge in different languages usually acti-
075 vates a shared set of neurons in feed-forward net-
076 works (FFNs), which we call language-agnostic
077 factual neurons. These neurons represent the se-
078 mantic correlations among the same multilingual
079 factual knowledge and are located in certain lay-
080 ers. Inspired by this finding, we propose a new
081 MKE method by locating and modifying Language-
082 Agnostic Factual Neurons (LAFN) to edit multi-
083 lingual knowledge simultaneously. Specifically,
084 we generate a set of paraphrases for each multilin-
085 gual knowledge to be edited to precisely locate the
086 corresponding language-agnostic factual neurons.
087 Then we optimize the update values for modify-
088 ing these located neurons to achieve simultaneous
089 modification of the same multilingual knowledge.
090 Additionally, to avoid the degradation of the edited
091 model’s general abilities due to directly modify-
092 ing model parameters (Gu et al., 2024), we do not
093 update the model parameters but store the update
094 values of the edited neurons in the cache. When the
095 edited subject appears in the user query, the relative
096 update values will be retrieved and used for model
097 inference.

098 To evaluate the effectiveness of our method, we
099 conduct experiments on two multilingual bench-
100 marks, Bi-ZsRE (Wang et al., 2023a) and MzsRE
101 (Wang et al., 2023b). Experimental results demon-
102 strate that our method outperforms existing MKE
103 methods in terms of Reliability, Generality, and Lo-
104 cality, indicating the importance of considering the
105 inner semantic connections between multilingual
106 knowledge.

107 In summary, the major contributions of this pa-
108 per are as follows¹:

- 109 • We propose a new MKE method by locating
110 and modifying language-agnostic factual neu-
111 rons that represent the deep semantic connec-
112 tions between multilingual knowledge.
- 113 • Experimental results on Bi-ZsRE and MzsRE
114 benchmarks demonstrate that our method
115 achieves outstanding edit performance, indi-
116 cating the effectiveness of our method.

¹The code will be released after acceptance.

- We discover that the language-agnostic fac-
tual neurons in the middle layers are crucial
for achieving MKE, shedding light on com-
prehension of the multilingual capabilities of
LLMs.

2 Methodology

In this section, we first give the definition of MKE (§2.1). Then we investigate how LLMs handle factual knowledge of different languages by identifying and analyzing the associated neurons (§2.2). Subsequently, we introduce our method LAFN for multilingual knowledge editing (§2.3).

2.1 Task Definition

MKE aims to simultaneously update multilingual knowledge with new information while preserving previous accurate knowledge within the model. Formally, we denote the original model as \mathcal{F}_θ and the multilingual group of an edit descriptor (x^e, y^e) as $G = \{\ell \in L | (x_\ell^e, y_\ell^e)\}$, where x_ℓ^e is the question for the knowledge to be edited in language ℓ and usually contains a subject and a relation, and y_ℓ^e is the new answer of x_ℓ^e . On this basis, MKE will lead to a model \mathcal{F}'_θ to correctly answer the edited question x_ℓ^e in each language ℓ and meanwhile maintain the original prediction on other unedited questions:

$$\forall \ell \in L, \mathcal{F}'_\theta(x_\ell) = \begin{cases} y_\ell^e, & x_\ell \in I(x_\ell^e), \\ \mathcal{F}_\theta(x_\ell), & x_\ell \notin I(x_\ell^e), \end{cases} \quad (1)$$

where $I(x_\ell^e)$ denotes a broad set of inputs with the same semantics as x_ℓ^e (Wang et al., 2023a).

2.2 Language-Agnostic Factual Neurons

Existing research has proven that knowledge neurons within FFNs store language-specific knowledge (Tang et al., 2024) and language-independent knowledge (Chen et al., 2023). And manipulating the values of these neurons has the potential to change the model’s behaviors, *e.g.*, changing the language-specific neurons can influence the language of the model’s output (Tang et al., 2024). Inspired by these findings, we first identify neurons associated with multilingual factual knowledge in two multilingual LLMs. Specifically, we separately identify the factual neurons for each language and then take the intersection of neurons for multiple languages as the language-agnostic factual neurons.

Identifying Language-Agnostic Factual Neurons. For most current LLMs (e.g., LLaMA2 (Touvron et al., 2023), Qwen (Bai et al., 2023), and Gemma (Team et al., 2024)), the calculation process of the i -th FFN layer can be formally described as:

$$h^i = (\text{act_fn}(\tilde{h}^i W_1^i) \otimes \tilde{h}^i W_2^i) \cdot W_3^i, \quad (2)$$

where \tilde{h}^i/h^i are the output hidden states of the i -th attention/FFN layer, $\text{act_fn}(\cdot)$ is the activation function, and W_1^i, W_2^i, W_3^i are the gate_proj, up_proj, down_proj matrix, respectively. In this process, knowledge neurons usually refer to the activations calculated by the activation function after the first matrix of FFNs, e.g., $\text{act_fn}(\tilde{h}^i W_1^i)$. Then we define that the j -th neuron in the i -th FFN layer is activated when $\text{act_fn}(\tilde{h}^i W_1^i)_j > 0$ following the previous work (Tang et al., 2024).

For the factual neurons of language ℓ , we use a factual corpus C_ℓ in language ℓ to track the activation of neurons in each FFN layer during the forward propagation. Subsequently, we identify and select the neurons that are activated most frequently to form the final neuron set. For instance, the set of factual neurons in the i -th FFN layer D_ℓ^i can be identified using C_ℓ as follows:

$$N^i = \{n_j^i | n_j^i = \sum_{c \in C_\ell} \mathbb{1}(\text{act_fn}(\tilde{h}_c^i W_1^i)_j > 0)\}, \quad (3)$$

$$D_\ell^i = \{j | \frac{n_j^i}{\max(N^i)} > \beta\}, \quad (4)$$

where \tilde{h}_c^i contains \tilde{h}^i at each token position in sentence c , $\mathbb{1}(\text{act_fn}(\tilde{h}_c^i W_1^i)_j > 0)$ equals to 1 when $\text{act_fn}(\tilde{h}_c^i W_1^i)_j > 0$ otherwise 0, n_j^i is the total activation counts of the j -th neuron in the i -th FFN layer, N^i is the set of activation counts of all neurons in i -th FFN layer when processing C_ℓ , and β is the threshold to control the amount of D_ℓ^i . After obtaining the sets of factual neurons for each language in L , we calculate the intersection of all these sets in the i -th FFN layer to extract the shared knowledge among all languages as follows:

$$D^i = D_{\ell_1}^i \cap D_{\ell_2}^i \cap \dots \cap D_{\ell_L}^i, \quad (5)$$

where we call D^i as the language-agnostic factual neurons in the i -th layer, implying the semantic connections of multilingual knowledge.

Experiments. We conduct analysis on PARAREL (Elazar et al., 2021), which contains factual knowledge with 34 relations in English. Here, we identify

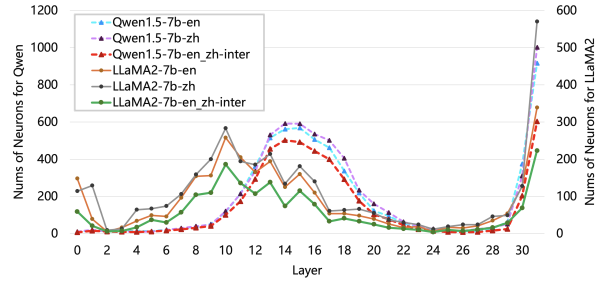


Figure 2: The identified neuron numbers in each layer of Qwen1.5-7b and LLaMA2-7b. “xxx-en” and “xxx-zh” represent the English and Chinese factual neurons respectively. “xxx-inter” refers to the language-agnostic factual neurons shared by English and Chinese.

the language-agnostic factual knowledge between English (*en*) and Chinese (*zh*). Firstly, we randomly choose 3000 sentences in each relation from PARAREL to build the factual corpus C_{en} (around 100k), and then utilize the Google Translate API to translate C_{en} to C_{zh} . We select two public multilingual LLMs: LLaMA2-7b (Touvron et al., 2023) and Qwen1.5-7b (Bai et al., 2023). The layer numbers of the two models are both 32. The threshold β in Eq.(4) is set to 0.8. According to Eq.(4) and Eq.(5), we count the language-agnostic factual neurons in each layer for the two LLMs.

Results. We plot the identified neuron numbers in each layer of the two models in Figure 2, including the factual neurons of each language and the language-agnostic factual neurons. It shows that the changes of the neuron numbers for the two models exhibit similar trends, with a greater presence of language-agnostic knowledge neurons in the middle layers and the last layer (refer to the green and red lines in Figure 2). The difference is that LLaMA2-7b peaks in quantity at the 10th layer, while Qwen1.5-7b reaches its peak at the 14th layer. And Qwen1.5-7b has more language-agnostic factual neurons than LLaMA2-7b. In conclusion, the experimental results prove the existence of language-agnostic factual neurons, which represent the deep semantic connections between the same factual knowledge in different languages and are mainly located in certain layers. Based on this finding, we design a method by locating and modifying language-agnostic factual neurons to edit multilingual knowledge simultaneously.

2.3 LAFN

Figure 3 shows the architecture of our method. We first locate the language-agnostic factual neurons

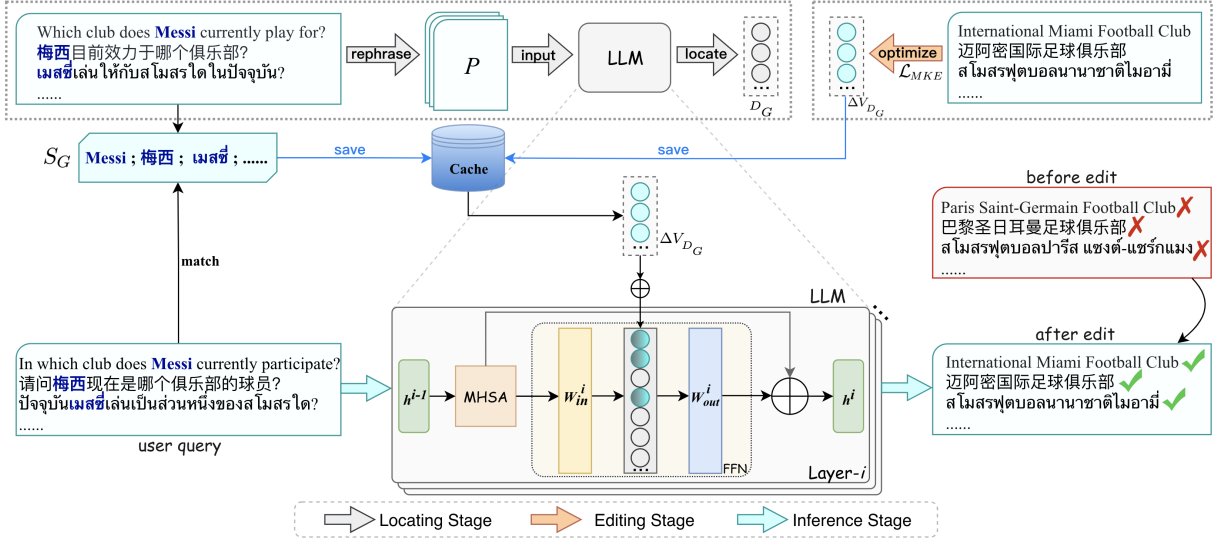


Figure 3: The architecture of LAFN. Given the multilingual knowledge to be edited (including the aligned multilingual subject set S_G), we first locate the corresponding language-agnostic neurons D_G . Then the update values ΔV_{D_G} is optimized for modifying D_G , and $\{S_G : \Delta V_{D_G}\}$ is stored in cache. When the subject of the user query is matched in the cache, the relative ΔV_{D_G} is used for model inference.

for each group of multilingual edit descriptors, and then we optimize the update values to modify these neurons and store them in the cache. During the inference stage, when the subject of the user query is matched in the cache, the relative update values are utilized for model inference.

During the locating stage, given the multilingual group G of an edit descriptor (x^e, y^e) ($G = \{\ell \in L | (x_\ell^e, y_\ell^e)\}$), we first locate the factual neurons D_ℓ^i in i -th layer for (x_ℓ^e, y_ℓ^e) in language ℓ according to Eq.(3) and Eq.(4). Specifically, to more precisely locate the neurons that are semantically related to x_ℓ^e , we use an LLM to generate several paraphrases for x_ℓ^e to build its paraphrase set as the factual corpus C_ℓ in Eq.(3). After obtaining D_ℓ^i in each language ℓ , we calculate the language-agnostic factual neuron set D_G^i of G in i -th layer following Eq.(5).

During the editing stage, given one multilingual edit description group G and its located language-agnostic factual neuron set D_G , we aim to modify the values of D_G to edit knowledge in G simultaneously. Following the settings of MEMIT (Meng et al., 2023) and PMET (Li et al., 2024), we modify the values V_{D_G} of D_G at the last token position of the subject in the question x_ℓ^e . As for subjects, we obtain the corresponding aligned multilingual subject set S_G from G (refer to S_G in Figure 3). Then we will optimize the update values ΔV_{D_G} for adding to V_{D_G} to achieve MKE. That is, the model should generate the corresponding new answer y_ℓ^e

by adding the ΔV_{D_G} :

$$\mathcal{F}_{(\theta, V_{D_G} + \Delta V_{D_G})}(x_\ell^e) = y_\ell^e \quad (6)$$

To this end, we calculate the \mathcal{L}_{target} to optimize ΔV_{D_G} :

$$\mathcal{L}_{target} = \frac{1}{|L|M} \sum_{\ell \in L} \sum_{m=1}^M -\log P_{\mathcal{F}_\theta}(y_\ell^e | p_\ell^m + x_\ell^e), \quad (7)$$

where $\ell \in L$, $\mathcal{F}_\theta = \mathcal{F}_{(\theta, V_{D_G} + \Delta V_{D_G})}$, and p_ℓ^m represents a randomly generated prefix to improve generalization (Meng et al., 2023) on $I(x_\ell^e)$, and M is the total number of prefixes.

Additionally, to ensure that the knowledge under the other relations of S_G is not affected, we also use \mathcal{L}_{kl} to optimize ΔV_{D_G} similar to MEMIT (Meng et al., 2023) and PMET (Li et al., 2024):

$$\mathcal{L}_{kl} = \frac{1}{|L|} \sum_{\ell \in L} \text{KL}[P_{\mathcal{F}_\theta}(y | q_\ell) || P_{\mathcal{F}_\theta'}(y | q_\ell)], \quad (8)$$

where q_ℓ has the format of “ $\{s_\ell\}$ is a” in language ℓ , s_ℓ is the subject in x_ℓ^e and $s_\ell \in S_G$, and $\text{KL}[\cdot || \cdot]$ is the Kullback-Leibler divergence (Kullback and Leibler, 1951).

In the end, the overall optimized objective \mathcal{L}_{MKE} consists of the above two loss functions:

$$\mathcal{L}_{MKE} = \lambda_1 \mathcal{L}_{target} + \lambda_2 \mathcal{L}_{kl}, \quad (9)$$

where λ_1 and λ_2 are hyperparameters to control the weight of two loss functions.

After obtaining ΔV_{DG} , we store $\{S_G : \Delta V_{DG}\}$ in the cache to avoid directly modifying the model parameters. When the subject s_ℓ of the current query x_ℓ is matched² in S_G , we retrieve the corresponding ΔV_{DG} for model inference as follows:

$$\mathcal{F}'_\theta(x_\ell) = \begin{cases} \mathcal{F}_{(\theta, V_{DG} + \Delta V_{DG})}(x_\ell), & s_\ell \in S_G. \\ \mathcal{F}_\theta(x_\ell), & s_\ell \notin S_G. \end{cases} \quad (10)$$

3 Experiments

3.1 Experimental Settings

Datasets and Metrics. We conduct our experiments on Bi-ZsRE (Wang et al., 2023a) and MzsRE (Wang et al., 2023b). Bi-ZsRE covers English (*en*) and Chinese (*zh*) languages, and each language contains 10000/3000/1037 samples for the train/dev/test set. MzsRE covers 12 languages: English (*en*), Chinese (*zh*), Czech (*cz*), German (*de*), Dutch (*nl*), Spanish (*es*), French (*fr*), Portuguese (*pt*), Russian (*ru*), Thai (*th*), Turkish (*tr*), and Vietnamese (*vi*). And each language consists of 10000/743 examples for the train/test set. Following Wang et al. (2023a), we calculate the F1 value of Reliability, Generality, Locality, and Portability as our evaluation metrics.

Backbones. In our experiments, we select two strong multilingual models LLaMA2-7b (Touvron et al., 2023) and Qwen1.5-7b (Bai et al., 2023) as backbones to conduct MKE. LLaMA2-7b is a widely used backbone known for its excellent universal capabilities. Qwen1.5-7b exhibits a strong foundational capability and demonstrates superior performance specifically in Chinese³.

Implementation Details. When locating neurons in §2.3, we utilize the Qwen1.5-14b-Chat⁴ model to generate 30 paraphrases for x_ℓ^e . The detailed instruction is listed in Appendix A. The threshold β in Eq.(4) is set to 0.1. The length of each randomly generated prefix p_ℓ^m in Eq.(7) is set to 5, and the total amount M of prefixes for each language is set to 4. Additionally, λ_1 is set to 1, and λ_2 is set to 0.0625. We use the Adam optimizer (Kingma and Ba, 2017) with a learning rate of 5e-1 during training. For layers to be modified, we set (10, 11, 12) for LLaMA2-7b and (14, 15, 16) for Qwen1.5-7b, respectively.

²Here, we use the exact-match method.

³<https://qwenlm.github.io/zh/blog/qwen1.5/>

⁴<https://huggingface.co/Qwen/Qwen1.5-14B-Chat>

3.2 Contrast Methods

Fune-tuning Method. We directly use LoRA (Hu et al., 2021) to conduct parameter-efficient tuning for the original model, namely LoRA-FT.

MKE Method⁵. ReMaKE (Wang et al., 2023b) retrieves related knowledge from a multilingual knowledge base as the context to instruct the model. Here, for the language to be tested, we separately retrieve one question with the answer from each other language as the context.

Adaptations of KE methods. We mainly adapt some Locate-then-Edit methods to MKE. For example, ROME (Meng et al., 2022) modifies the output matrix of one FFN layer located following causal tracing analysis. MEMIT (Meng et al., 2023) updates the output matrix of multiple layers simultaneously for supporting batch editing. PMET (Li et al., 2024) conducts more precise editing based on MEMIT. We extend ROME, MEMIT, and PMET to M-ROME, M-MEMIT, and M-PMET to edit multilingual knowledge simultaneously. Specifically, since the knowledge to be edited of different languages corresponds to different answers, we train the new value for updating FFNs separately for each language. And we estimate the previously memorized keys of FFNs for each language.

3.3 Experimental Results

Results on Bi-ZsRE. Table 1 shows the results on Bi-ZsRE using LLaMA2-7b and Qwen1.5-7b as backbones. From the “*avg*” column, the average results of all metrics demonstrate that our method outperforms other baselines significantly, indicating the importance of considering the deep semantic connections between multilingual knowledge. In terms of Reliability and Generality, our method exceeds other methods to a large extent. This superiority indicates that updating the language-agnostic factual neurons can edit the multilingual knowledge (needs to be edited) more effectively and generalize better on the equivalent questions that have the same semantics as the edited questions. LoRA-FT and ReMaKe perform poorly, while M-ROME, M-MEMIT, and M-PMET perform moderately among all methods. Specifically, M-ROME is less effective than M-MEMIT and M-PMET because it only updates a single layer. M-MEMIT and M-PMET have similar performances but are both inferior to

⁵The code of MPN is not open-source, and LiME is based on mBERT without exploring the generation task, so we do not compare these two methods.

Methods	Test on en				Test on zh				avg
	Reliability	Generality	Locality	Portability	Reliability	Generality	Locality	Portability	
LLaMA2-7b (Edit on en & zh)									
LoRA-FT	21.90	21.15	81.90	27.07	15.30	15.43	75.02	13.05	33.85
ReMaKe	32.90	33.78	100.00	28.35	31.78	31.66	99.94	15.77	46.77
M-ROME	69.48	64.42	96.19	26.27	37.94	35.61	91.41	10.38	53.96
M-MEMIT	84.73	74.13	98.70	28.65	41.58	38.18	97.63	11.38	59.37
M-PMET	85.40	77.02	98.31	29.30	41.25	37.80	97.60	10.88	59.70
LAFN (Ours)	98.66	93.80	100.00	30.93	56.22	53.42	100.00	12.72	68.22
Qwen1.5-7b (Edit on en & zh)									
LoRA-FT	20.31	20.50	84.04	24.91	32.95	32.59	88.38	33.53	42.15
ReMaKe	46.20	46.41	100.00	29.79	66.04	67.07	100.00	43.98	62.44
M-ROME	88.37	77.05	95.66	31.02	93.68	86.01	95.36	37.99	75.64
M-MEMIT	94.36	88.27	95.72	31.13	96.80	92.96	96.63	37.03	79.11
M-PMET	95.59	88.46	95.39	30.66	96.66	93.37	96.12	37.97	79.28
LAFN (Ours)	99.27	94.13	100.00	28.20	99.86	95.08	100.00	36.16	81.59

Table 1: The F1 results on Bi-ZsRE using LLaMA2-7b and Qwen1.5-7b as backbones. Results highlighted in **bold** represent the best results. “*avg*” denotes the average value of all metrics in both two languages.

our method, demonstrating that the simple adaptations of these methods to MKE are less effective. As for Locality, both our method and REMAKE achieve the “100.00” value since the two methods do not modify the parameters of the original model during the editing process, not influencing previously learned knowledge. While the other methods modify the model parameters and result in lower Locality scores. Among them, LoRA-FT dramatically modifies the model, scoring the lowest.

Portability, as a more difficult metric, measures whether the edited model can reason based on the edited knowledge via a portability question (Yao et al., 2023). The corresponding results show that all methods underperform on this metric. Our method achieves the best result on the English test set when editing LLaMA2-7b, and M-MEMIT performs best on the English test set when editing Qwen1.5-7b. ReMaKe achieves the best results on the Chinese test set since the longer context improves the reasoning ability of LLMs. However, there is still substantial room for all methods to enhance the reasoning ability based on edited knowledge. Moreover, we observe that Qwen1.5-7b exhibits notably superior edit performance in Chinese compared to LLaMA2-7b, indicating that the inherent language capabilities of a model have a crucial impact on its edit performance.

Results on MzsRE. As for the more challenging scenarios, the average results of 12 languages on MzsRE are reported in Table 2 (using LLaMA2-7b as the backbone). The results show that our method obtains the best overall performance, proving the effectiveness of updating the language-agnostic fac-

Methods	Reliability	Generality	Locality	Portability	avg
LoRA-FT	24.03	23.94	64.74	22.64	33.84
ReMaKe	41.86	42.37	100.00	26.36	52.65
M-ROME	32.96	32.20	62.40	11.94	34.87
M-MEMIT	76.51	70.24	93.26	23.14	65.79
M-PMET	72.79	69.10	93.32	22.51	64.43
LAFN (Ours)	85.79	80.75	100.00	22.47	72.25

Table 2: The average F1 results of 12 languages on MzsRE using LLaMA2-7b as the backbone. Results highlighted in **bold** represent the best results. “*avg*” denotes the average value of all metrics in 12 languages.

tual neurons. Specifically, LAFN surpasses other methods in terms of Reliability, Generality, and Locality by a large margin. Additionally, “M-ROME” performs much worse in 12 languages than in just two languages, demonstrating that this method struggles to support simultaneous editing of more language knowledge due to the limited edit region. Detailed results of each language are listed in Table 6 of Appendix B.

4 Analysis

In §4.1, we initially analyze the performance under different layer settings. Then we compare different locating strategies to prove that using paraphrases during the locating stage can improve the edit performance (§4.2). Subsequently, we investigate the impact of our method on the unedited knowledge of the edited subjects (§4.3).

4.1 Different Layer Settings

In this section, we explore how editing performance changes when editing different layers. Figure 2 in §2.2 shows that the language-agnostic factual neu-

A Single Layer	avg	Multiple Layers	avg
0	42.23	2-10	65.73
2	57.91	10-31	65.53
10	65.38	10-11	67.81
13	64.95	10-11-12	68.22
24	57.73	10-11-12-13	67.92
31	36.98	10-11-12-13-31	67.67

Table 3: The results of different layer settings on Bi-ZsRE using LLaMA2-7b as the backbone.

441 rons are mostly in some middle layers and the last
442 layer of all FFNs. To investigate the correlation
443 between edit performance and edited layers, we
444 conduct our method in different layer settings ac-
445 cording to the number of language-agnostic factual
446 neurons, including a single layer and multiple lay-
447 ers. The corresponding results reported in Table 3
448 show that in the single-layer setting, the edit per-
449 formance achieves best in the 10th layer and worst
450 in the last layer (31th). Although the last layer also
451 has numerous language-agnostic factual neurons,
452 we conjecture that these neurons are directly re-
453 lated to the final outputs, and thus a single update
454 vector is difficult to fulfill answers in all languages.
455 Moreover, we simultaneously edit multiple layers
456 based on the 10th layer, and the results show that
457 editing multiple layers can further improve edit
458 performance, with the best performance observed
459 in (10, 11, 12) layers. In short, these results sug-
460 gest that language-agnostic factual neurons in the
461 middle layers are crucial for achieving MKE.

4.2 Different Locating Strategies

462 To verify the effectiveness of using paraphrases
463 during the locating stage, we compare three differ-
464 ent locating strategies with the original LAFN: (1)
465 (no-PGs) not using paraphrases to assist in locat-
466 ing neurons, *i.e.*, only using a single sentence to
467 locate neurons; (2) (all) modifying all neurons of
468 the same layers as LAFN without locating concrete
469 knowledge-related neurons; (3) (random) randomly
470 selecting the same number of neurons in the same
471 layers as LAFN to modify. The results listed in
472 Table 4 show that the performance of the three
473 settings declines compared to the original LAFN,
474 particularly regarding Generality and Portability.
475 Although the results of Reliability with “no-PGs”
476 and “all” have a slight improvement, the results of
477 Generality and Portability decline obviously due
478 to the modified neurons being too limited or too
479 broad. In the “random” setting, the results of all
480

Methods	Reliability	Generality	Portability	avg
LAFN (Ours)	77.44	73.61	21.83	68.22
(no-PGs)	77.61 ↑	73.35 ↓	21.54 ↓	68.12 ↓
(all)	77.47 ↑	73.55 ↓	21.69 ↓	68.18 ↓
(random)	77.42 ↓	69.99 ↓	21.16 ↓	67.14 ↓

Table 4: The results of different locating strategies on Bi-zsRE using LLaMA2-7b as the backbone. The **Locality** scores are all 100 for these settings and thus not listed. “**avg**” averages the scores of these 4 metrics.

Methods	Test on en	Test on zh	avg
M-ROME	92.91	96.50	94.71
M-MEMIT	94.23	97.33	95.78
M-PMET	93.81	97.07	95.44
LAFN (Ours)	94.80	98.19	96.50

Table 5: The F1 scores of Locality-Hard based on Bi-zsRE using LLaMA2-7b as the backbone.

481 metrics have notably decreased compared to the
482 original LAFN. To sum up, these results prove that
483 using paraphrases during the locating stage can
484 enhance the located neurons more semantically rel-
485 evant to the multilingual knowledge to be edited,
486 thus improving the edit performance.

4.3 Impact on Unedited Knowledge of the Edited Subjects

487 **Locality-Hard.** During inference, we directly add
488 the corresponding update values to the last token
489 position of the subject when the current subject is
490 matched in the cache. This process may have a
491 side effect on the unedited knowledge related to
492 the edited subjects (*e.g.*, knowledge with the same
493 subject as the edited knowledge but different re-
494 lations). Therefore, we investigate whether our
495 method harms this type of knowledge. Specifically,
496 we collect some extra knowledge to build a more
497 challenging test set, which has the same subjects as
498 each edited example but different relations (please
499 refer to details in Appendix C). Then we calculate
500 the Locality metric on this test set and denote it as
501 Locality-Hard. The results in Table 7 show that our
502 method achieves the highest score of Locality-Hard
503 compared to other methods. These results reflect
504 less impact of our method on the unedited knowl-
505 edge of the edited subjects, and also indicate that
506 the modified neurons by our method are strongly
507 related to the edited knowledge.
508

509 **Case Study.** To investigate the language-agnostic
510 factual neurons more clearly, we visualize the dif-
511 ferences between the neurons located by differ-
512 ent knowledge in the 10th FFN layer (as shown
513

in Figure 4). We list the selected knowledge in Table 7 of Appendix D. Specifically, we first locate the set of language-agnostic factual neurons D in the 10th FFN layer for each instance according to Eq.(5) and calculate the difference dif between two sets D_a and D_b following $dif = 1 - (\frac{|D_a \cap D_b|}{|D_a|} + \frac{|D_a \cap D_b|}{|D_b|})/2$. And $dif = 0$ represents $D_a = D_b$, with the darkest color in Figure 4. $dif = 1$ represents $D_a \cap D_b = \emptyset$ and the corresponding color is lightest. From Figure 4, we can observe several phenomena: (1) Instances of the same subject with the same relation have the smallest differences between their located neurons, which have the darkest color (refer to the region of the orange box). (2) Instances of the same subject but different relations have a small degree of differences between the located neurons, and the color is also relatively dark (refer to the region of the blue box). (3) Instances of different subjects have large differences between the located neurons, and the color is much lighter (refer to the region of the pink box). In summary, these differences in Figure 4 indicate that the neurons modified by our method are highly associated with the edited knowledge, bringing less impact on other knowledge.

5 Related Work

Multilingual Knowledge Editing. Existing MKE methods mostly adapt the monolingual KE methods to multilingual scenarios, overlooking the connections of multilingual knowledge. For example, LiME (Xu et al., 2023) proposes an editing framework using the parallel corpus to train hyper-networks, adapting the monolingual meta-learning edit methods to the cross-lingual scenario, such as KE (De Cao et al., 2021) and MEND (Mitchell et al., 2022). ReMaKE (Wang et al., 2023b) retrieves the multilingual aligned knowledge from a multilingual knowledge base as context to achieve MKE. Additionally, MPN (Si et al., 2024) trains multilingual patch neurons to store multilingual knowledge following T-Patcher (Huang et al., 2023), which only applies to classification tasks. By contrast, our method first locates the language-agnostic factual neurons using the knowledge to be edited and then modifies them, which considers the deep connections of multilingual knowledge and is more intuitive.

Multilingual Knowledge Analysis. Analyzing the multilingual capabilities of language models is always a research hotspot (Pires et al., 2019;

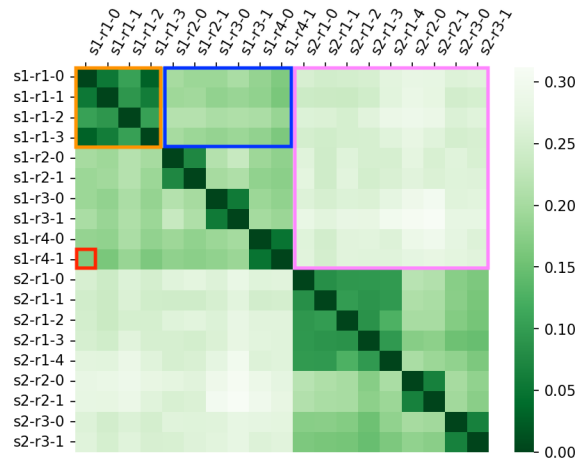


Figure 4: The differences between the language-agnostic factual neurons located by different knowledge in the 10th FFN layer. “s1” and “s2” represent two subject groups. “r1/2/3/4” are different relations under each subject. Each small square (refer to the red box) represents the difference dif between the two neuron sets, and the darker the color, the smaller the difference between the two sets.

Chai et al., 2022; Bhattacharya and Bojar, 2023; Kojima et al., 2024; Zhao et al., 2024), especially exploring the relationship between model architecture and multilingual capabilities. Tang et al. (2024) indicate that LLMs’ proficiency in processing a particular language is predominantly due to a small subset of neurons. Similar to our work, Chen et al. (2023) discover the language-independent knowledge neurons of mBERT and mGPT, which store knowledge in a form that transcends language, but ignores how to control neurons to achieve desired outputs. Differently, we first investigate the language-agnostic knowledge neurons related to specific fact knowledge in LLMs and then modify them to achieve multilingual knowledge editing.

6 Conclusion

In this work, we propose a new method LAFN to conduct multilingual knowledge editing by locating and modifying language-agnostic factual neurons. The experimental results on two benchmarks demonstrate our method outperforms existing MKE methods, indicating the effectiveness of our method and the importance of considering the semantic connections between multilingual knowledge. Furthermore, we find that language-agnostic factual neurons in the middle layers are crucial for MKE, which can provide insights into understanding the multilingual capabilities of LLMs.

592 Limitations

593 In our approach, it is necessary to provide the
594 aligned multilingual knowledge to be edited and
595 their corresponding multilingual subjects, which
596 is directly available in both Bi-ZsRE and MzsRE
597 datasets. However, for other datasets that do not
598 contain this information, we first need to preprocess
599 the data to support our method. For example,
600 if there is no corresponding multilingual data
601 available, using translation API can translate the
602 existing knowledge to be edited to other languages.
603 If the corresponding subjects are not annotated, existing
604 LLMs can be utilized to identify the aligned
605 multilingual subjects in the sentences of each language.
606 These preprocessing steps can be easily implemented
607 by calling existing tools. Moreover, the current method
608 for determining whether a subject exists in the cache
609 adopts the exact-match approach, which is too strict.
610 We will optimize it to a fuzzy matching method in
611 future work to enhance the performance in practical
612 application scenarios.

613 References

614 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
615 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
616 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
617 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
618 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
619 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
620 Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-
621 guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang,
622 Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,
623 Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan
624 Zhang, Yichang Zhang, Zhenru Zhang, Chang
625 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang
626 Zhu. 2023. Qwen technical report. *arXiv preprint*
627 *arXiv:2309.16609*.

628 Himanshu Beniwal, Kowsik D, and Mayank Singh.
629 2024. *Cross-lingual editing in multilingual language*
630 *models*. In *Findings of the Association for Computational*
631 *Linguistics: EACL 2024*, pages 2078–2128,
632 St. Julian’s, Malta. Association for Computational
633 Linguistics.

634 Sunit Bhattacharya and Ondřej Bojar. 2023. *Unveil-*
635 *ing multilinguality in transformer models: Exploring*
636 *language specificity in feed-forward networks*. In
637 *Proceedings of the 6th BlackboxNLP Workshop: Analyzing*
638 *and Interpreting Neural Networks for NLP*,
639 pages 120–126, Singapore. Association for Computational
640 Linguistics.

641 Yuan Chai, Yaobo Liang, and Nan Duan. 2022. *Cross-*
642 *lingual ability of multilingual masked language mod-*
643 *els: A study of language structure*. In *Proceedings*
644 *of the 60th Annual Meeting of the Association for*

Computational Linguistics (Volume 1: Long Papers),
645 pages 4702–4712, Dublin, Ireland. Association for
646 Computational Linguistics. 647

648 Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and
649 Jun Zhao. 2023. *Journey to the center of the knowl-*
650 *edge neurons: Discoveries of language-independent*
651 *knowledge neurons and degenerate knowledge neu-*
652 *rons*. *Preprint*, arXiv:2308.13198.

653 Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. *Edit-*
654 *ing factual knowledge in language models*. In *Pro-*
655 *ceedings of the 2021 Conference on Empirical Meth-*
656 *ods in Natural Language Processing*, pages 6491–
657 6506, Online and Punta Cana, Dominican Republic.
658 Association for Computational Linguistics.

659 Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha
660 Ravichander, Eduard Hovy, Hinrich Schütze, and
661 Yoav Goldberg. 2021. *Measuring and improving*
662 *consistency in pretrained language models*. *Preprint*,
663 arXiv:2102.01017.

664 Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-
665 Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024.
666 *Model editing can hurt general abilities of large lan-*
667 *guage models*. *Preprint*, arXiv:2401.04700.

668 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
669 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
670 Weizhu Chen. 2021. *Lora: Low-rank adaptation of*
671 *large language models*. *Preprint*, arXiv:2106.09685.

672 Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou,
673 Wenge Rong, and Zhang Xiong. 2023. *Transformer-*
674 *patcher: One mistake worth one neuron*. In *The*
675 *Eleventh International Conference on Learning Rep-*
676 *resentations*.

677 Diederik P. Kingma and Jimmy Ba. 2017. *Adam:*
678 *A method for stochastic optimization*. *Preprint*,
679 arXiv:1412.6980.

680 Takeshi Kojima, Itsuki Okimura, Yusuke Iwasawa, Hit-
681 omi Yanaka, and Yutaka Matsuo. 2024. *On the multi-*
682 *lingual ability of decoder-based pre-trained language*
683 *models: Finding and controlling language-specific*
684 *neurons*. *Preprint*, arXiv:2404.02431.

685 S. Kullback and R. A. Leibler. 1951. *On Information*
686 *and Sufficiency*. *The Annals of Mathematical Statis-*
687 *tics*, 22(1):79 – 86.

688 Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun
689 Ma, and Jie Yu. 2024. *Pmet: Precise model editing*
690 *in a transformer*. *Preprint*, arXiv:2308.08742.

691 Kevin Meng, David Bau, Alex Andonian, and Yonatan
692 Belinkov. 2022. *Locating and editing factual asso-*
693 *ciations in gpt*. In *Advances in Neural Information*
694 *Processing Systems*, volume 35, pages 17359–17372.
695 Curran Associates, Inc.

696 Kevin Meng, Arnab Sen Sharma, Alex J Andonian,
697 Yonatan Belinkov, and David Bau. 2023. *Mass-*
698 *editing memory in a transformer*. In *The Eleventh*

699					
700		<i>International Conference on Learning Representations.</i>			
701	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2022. Fast model editing at scale. In <i>International Conference on Learning Representations.</i>				
702					
703					
704					
705	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001, Florence, Italy. Association for Computational Linguistics.				
706					
707					
708					
709					
710					
711	Nianwen Si, Hao Zhang, and Weiqiang Zhang. 2024. Mpn: Leveraging multilingual patch neuron for cross-lingual model editing. <i>Preprint</i> , arXiv:2401.03190.				
712					
713					
714	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. Language-specific neurons: The key to multilingual capabilities in large language models. <i>Preprint</i> , arXiv:2402.16438.				
715					
716					
717					
718					
719	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology. <i>Preprint</i> , arXiv:2403.08295.				
720					
721					
722					
723					
724					
725					
726					
727					
728					
729					
730					
731					
732					
733					
734					
735					
736					
737					
738					
739					
740					
741					
742					
743					
744					
745					
746					
747					
748					
749					
750					
751					
752					
753					
754					
755					
756	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, 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. Llama 2: Open foundation and fine-tuned chat models. <i>Preprint</i> , arXiv:2307.09288.				
757					
	Jiaan Wang, Yunlong Liang, Zengkui Sun, Yuxuan Cao, and Jiarong Xu. 2023a. Cross-lingual knowledge editing in large language models. <i>Preprint</i> , arXiv:2309.08952.				
	Weixuan Wang, Barry Haddow, and Alexandra Birch. 2023b. Retrieval-augmented multilingual knowledge editing. <i>Preprint</i> , arXiv:2312.13040.				
	Yang Xu, Yutai Hou, Wanxiang Che, and Min Zhang. 2023. Language anisotropic cross-lingual model editing. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5554–5569, Toronto, Canada. 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.				
	Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. 2024. How do large language models handle multilingualism? <i>Preprint</i> , arXiv:2402.18815.				
	A The Instruction for generating paraphrases.				
	We utilize the Qwen1.5-14b-Chat model to generate the paraphrase set P_ℓ for more precisely locating neurons. The English version of the instruction for inputting Qwen1.5-14b-Chat is “You are an expert at sentence rewriting. Below I will give you a subject and a question containing the subject. Please give me 30 questions including this subject				

Methods	cz	vi	tr	fr	es	zh	en	de	ru	nl	pt	th	avg
Reliability													
LoRA-FT	24.08	28.28	23.22	23.39	22.82	16.71	20.25	22.22	28.75	23.85	24.34	30.42	24.03
ReMaKe	38.23	46.82	42.01	38.55	37.56	29.98	33.44	36.32	54.28	36.91	38.91	69.32	41.86
M-ROME	30.97	29.86	24.67	34.80	32.61	20.94	41.77	35.20	31.93	36.09	31.73	44.91	32.96
M-MEMIT	83.20	73.61	71.12	81.21	81.03	39.23	82.10	81.93	90.66	79.74	78.28	76.04	76.51
M-PMET	75.67	72.27	67.88	76.19	74.08	37.43	83.44	77.97	86.02	76.58	74.08	71.85	72.79
LAFN	92.16	88.00	89.85	90.80	90.39	48.59	91.46	90.82	90.24	90.75	90.33	76.13	85.79
Generality													
LoRA-FT	23.88	28.02	22.92	22.92	22.86	16.99	20.03	22.43	28.81	23.59	24.11	30.70	23.94
ReMaKe	39.21	47.25	42.81	38.63	37.97	31.01	33.50	37.15	55.70	37.67	39.04	68.48	42.37
M-ROME	30.79	30.52	25.56	34.36	32.51	19.40	41.46	35.37	30.76	34.94	30.22	40.47	32.20
M-MEMIT	75.55	68.26	67.38	75.67	75.50	36.38	74.47	75.27	83.76	73.66	71.48	65.47	70.24
M-PMET	72.09	69.39	65.88	73.63	71.72	35.79	79.66	74.38	81.00	72.85	69.77	63.07	69.10
LAFN	87.65	82.27	85.19	86.60	86.54	46.08	87.41	86.01	83.38	84.48	84.55	68.85	80.75
Locality													
LoRA-FT	64.08	62.35	57.47	63.01	74.12	64.23	80.84	68.09	62.77	61.69	59.91	58.27	64.74
ReMaKe	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
M-ROME	55.16	63.19	47.65	64.13	71.89	64.20	83.89	69.78	54.90	61.85	61.37	50.82	62.40
M-MEMIT	94.19	93.92	90.82	94.81	95.99	95.36	97.48	94.82	90.58	93.68	92.98	84.48	93.26
M-PMET	94.31	93.45	90.53	95.05	95.97	95.02	97.55	95.25	91.31	93.83	93.84	83.70	93.32
LAFN	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Portability													
LoRA-FT	22.40	28.87	19.56	22.82	21.2	11.49	20.94	21.80	29.25	22.70	21.76	28.92	22.64
ReMaKe	26.31	33.82	22.87	26.82	25.23	14.47	27.39	25.19	34.06	24.14	25.19	30.86	26.36
M-ROME	10.57	14.95	9.07	11.43	10.55	5.64	14.93	10.85	15.26	10.94	9.92	19.15	11.94
M-MEMIT	23.38	29.55	21.30	23.50	22.45	10.75	27.78	22.37	26.88	22.42	21.87	25.41	23.14
M-PMET	22.13	28.99	20.90	22.34	21.41	10.12	28.13	21.85	26.86	21.76	21.00	24.63	22.51
LAFN	22.30	28.84	20.93	22.78	22.06	10.53	27.45	22.14	25.19	20.83	21.22	25.32	22.47

Table 6: The F1 results on the MzsRE dataset using LLaMA2-7b as the backbone.

in English. They must have the same semantics as the given question. Subject: {}. Question containing this Subject: {}”.

B Detailed Results on MzsRE

The detailed results of each language on MzsRE are listed in Table 6.

C Details for Locality-Hard

To investigate whether our method harms the unedited knowledge of the edited subjects, we call Qwen-max API to collect some knowledge with the same subject as each edited example but different relations based on the test set of Bi-ZsRE. Notably, the Qwen-max API can use the searched results to enhance the accuracy of the generated answers. We use these collected questions to build the challenging test set. Then we calculate the Locality metric on this test set and denote it as Locality-Hard.

D Selected Cases

The selected English examples in Figure 4 are listed in Table 7. Since there is a one-to-one correspondence between Chinese and English examples, we do not list Chinese examples again.

s1-en	Alec Rose
s1-r1-0-en	What war did Alec Rose participate in?
s1-r1-1-en	In what war did Alec Rose fight?
s1-r1-2-en	What war or battle involved Alec Rose?
s1-r1-3-en	Which war was Alec Rose in?
s1-r2-0-en	Where was Alec Rose born?
s1-r2-1-en	Alec Rose was born in which location?
s1-r3-0-en	When did Alec Rose receive the MBE?
s1-r3-1-en	In what year was Alec Rose awarded the MBE?
s1-r4-0-en	What was Alec Rose's profession?
s1-r4-1-en	In what field was Alec Rose employed?
s2-en	Elk's Head of Huittinen
s2-r1-0-en	When was Elk's Head of Huittinen discovered?
s2-r1-1-en	When was the discovery of Elk's Head of Huittinen?
s2-r1-2-en	What year was Elk's Head of Huittinen discovered?
s2-r1-3-en	When did the discovery or creation of Elk's Head of Huittinen occur?
s2-r1-4-en	Could you provide the year when the landmark Elk's Head in Huittinen was first brought to light?
s2-r2-0-en	In which country is Elk's Head of Huittinen located?
s2-r2-1-en	To which nation does Elk's Head of Huittinen belong?
s2-r3-0-en	What is the historical significance of Elk's Head of Huittinen?
s2-r3-1-en	What role does Elk's Head of Huittinen play in the local history?

Table 7: The selected examples in English of Figure 4.