# Editing Knowledge Representation of Language Model via Rephrased Prefix Prompts

Yuchen Cai[1,2], Ding Cao[1,2], Rongxi Guo[1,2], Yaqin Wen[1,2],

Guiquan Liu[1,2(✉)], and Enhong Chen[1,2]

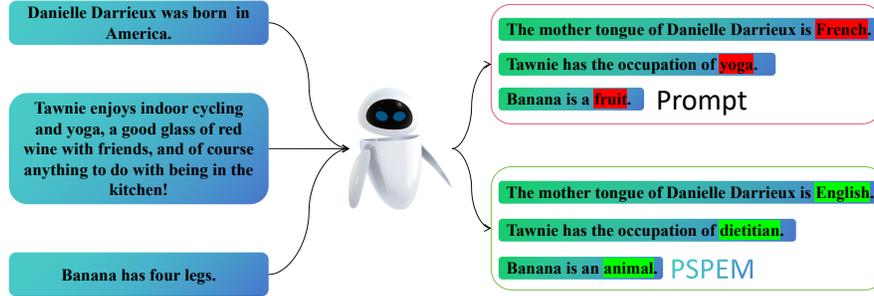[1] University of Science and Technology of China, Hefei, China
[2] State Key Laboratory of Cognitive Intelligence, Hefei, China
{caiyuchen, caoding, guorongxi, wyq_65}@mail.ustc.edu.cn,
{gqliu,cheneh}@ustc.edu.cn

**Abstract.** Neural language models (LMs) have been extensively trained on vast corpora to store factual knowledge about various aspects of the world described in texts. Current technologies typically employ knowledge editing methods or specific prompts to modify LM outputs. However, existing knowledge editing methods are costly and inefficient, struggling to produce appropriate text. Additionally, prompt engineering is opaque and requires significant effort to find suitable prompts. To address these issues, we introduce a new method called PSPEM (**P**refix **S**oft-**P**rompt **E**diting **M**ethod), that can be used for a lifetime with just one training. It resolves the inefficiencies and generalizability issues in knowledge editing methods and overcomes the opacity of prompt engineering by automatically seeking optimal soft prompts. Specifically, PSPEM adopts a prompt encoder and an encoding converter to compress and refine key information in prompts and adopts prompt alignment techniques to guide model generation, ensuring text consistency and adherence to the intended structure and content. We have validated the effectiveness of PSPEM through knowledge editing and attribute inserting. On the COUNTERFACT dataset, PSPEM achieved nearly 100% editing accuracy and demonstrated the highest level of fluency. We further analyzed the similarities between PSPEM and original prompts and their impact on the model's internals. The results indicate that PSPEM can serve as an alternative to original prompts, supporting the model in effective editing.

**Keywords:** Language model · Prompt learning · Knowledge editing · Knowledge representation

## 1 Introduction

Language models based on the Transformer architecture, such as GPT and BERT, have revolutionized various natural language processing tasks with their capacity to store and utilize real-world factual knowledge within their parameters [8, 24, 27]. For example, when asked, "Where is the Eiffel Tower located?"

**Fig. 1.** By inputting the prompts on the left side into the model, traditional prompt engineering generate erroneous text, while PSPEM can correct such errors.

GPT provides the accurate answer, "Paris". However, inconsistencies or biases present in the pre-training data can propagate into the text generated by the model, leading to errors or contradictions. Additionally, the knowledge of the world changes over time, which presents a challenge to the static knowledge in the model. Addressing this issue requires a nuanced approach to updating the model's knowledge base. Merely retraining the entire model with new data is prohibitively expensive and time-consuming, while fine-tuning, focusing solely on specific updated knowledge poses the risk of overfitting and compromising the model's ability to generalize.

Recent advancements propose a more dynamic and efficient approach to knowledge updating to mitigate these issues, called knowledge editing [29, 30]. This technology allows for selective updates and adjustments to the model's knowledge without retraining. These methods aim to balance the need for accurate, up-to-date information with the practical constraints of computational resources and time. The efficacy of knowledge editing is predominantly quantified by two pivotal metrics: generalization and specificity. Generalization entails the model's proficiency in extending the modified knowledge across a spectrum of analogous prompts, ensuring consistent application and understanding of the targeted information [7,32]. Conversely, specificity, also referred to as locality, necessitates the model's capacity to isolate the modification impact, safeguarding the unaltered knowledge from inadvertent alteration [21]. Several new benchmarks are attempting to assess the model's ability to reason with new knowledge, and they are extending these methodologies into the realms of knowledge graphs [4] and multimodal domains [3].
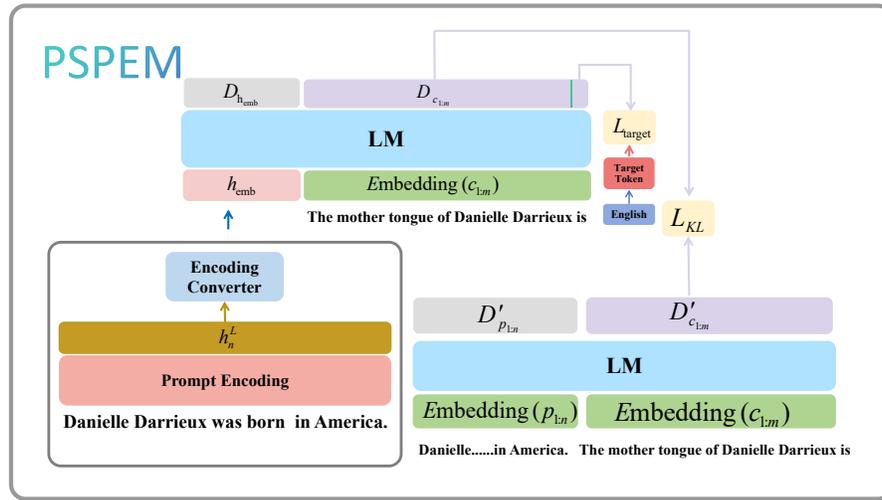
Methods of model editing fall into two distinct classifications depending on the alteration of the original model weights: weight-preserved and weight-modified methods [30]. Weight-preserved strategies typically necessitate the inclusion of extra content, while weight-modified techniques directly alter the model's weights. Weight modification methods include hypernetwork-based learning methods and direct optimization methods. Hypernetwork-based learning

methods, such as KE [7], MEND [21] and MALMEN [26], utilize a hypernetwork to predict essential updates to the model's weights. Although this technique is promising, it necessitates considerable computational investment for hypernetwork training and frequently diminishes in effectiveness with the increase in language model size [30]. Optimization method ROME [19] employs causal mediation analysis to identify the editing region and focuses on altering specific information via rank-one adjustments to individual matrices. MEMIT [20] adhered to a similar methodology, adeptly modifying several parameter matrices concurrently to facilitate the simultaneous alteration of 10,000 knowledge entities, and showcasing robust generalization and specificity. PMET [15] advanced this technique, refining MEMIT's capabilities for more precise editing. However, previous research indicates that minor modifications to the parameters of large language models can impact the model's ultimate behaviour [25], and these methods do not allow the model to use new knowledge for reasonable inference [13, 23].

Prompt engineering [16,17] enables the modification of models without necessitating extensive retraining. Altering input prompts, allows models to adapt to diverse tasks and domains, thereby conserving resources. Nonetheless, finding the most effective prompts typically requires considerable manual intervention and iterative experimentation, which can be time-consuming and inefficient. Moreover, discrepancies between the knowledge encapsulated in the prompts and the model's inherent knowledge can lead to erroneous or inconsistent outputs. As depicted in Figure 1, the language model (LM) is prompted with "Danielle Darrieux was born in America." Upon generating a continuation of this prompt, the LM erroneously asserts that Danielle Darrieux's native language is French, thereby contradicting the prior context. This error occurred because LM developed a memory for Danielle Darrieux's native language, French, during pretraining [1, 22].

In-context learning [2] is a paradigm that does not require retraining, where knowledge is acquired from directly connected demonstrations in the input context. Unlike traditional prompting engineering, the method is capable of learning contextual relationships from multiple given instances, enabling context-based model editing and providing an efficient, lightweight knowledge editing approach [31]. Although this method addresses the issue of inconsistent contextual information in the model, it requires searching for multiple guiding instances, which puts an additional burden on knowledge editing. REMEDI [13] injects domain-specific knowledge into language models by encoding factual prompts corresponding to knowledge attributes in the direction space. However, when factual prompts conflict with the knowledge already present within the model, the REMEDI method still struggles to handle such contradictions, resulting in inconsistencies or errors in the results.

To overcome the limitation of the weight-modified method in utilizing new knowledge for reasoning and the poor editing accuracy and laborious of the weight-preserved method, we proposed PSPEM (**P**refix **S**oft-**P**rompt **E**diting **M**ethod), an innovative strategy rooted in prompt engineering that can be used for a lifetime with just one training. This method allows for precise, nuanced

**Fig. 2.** Illustration of PSPEM. Given a knowledge prompt (Danielle Darrieux was born in America.) and continuation words (The mother tongue of Danielle Darrieux is), PSPEM constructed more accurately encoded information from the prompt to increase the probability of the target token (English).

modifications to the model's output by employing a single knowledge prompt, all without altering the model's parameters. Specifically, PSPEM utilizes a prompt encoder to extract information, an encoding converter to refine key information in prompts, and adopts prompt alignment techniques to guide model generation. It ensures the accuracy of the model's output by maximizing the probability assigned by the language model to the target token, and by aligning with the original prompt's influence on continuation words, it guarantees the fluency of the output text and a high degree of consistency between the generated text information and the prompt information. We conducted evaluations on two knowledge editing tasks and two attribute inserting tasks. In the tasks of knowledge editing, PSPEM achieved nearly 100% editing accuracy while ensuring the fluency and consistency of the generated text. In terms of attribute inserting, PSPEM reached the state-of-the-art. The model can make reasonable inferences using the given prompts and generate text that aligns with the prompt information. We then analyzed the parallels between PSPEM and original prompts, measuring their impact on model output from multiple perspectives. The experimental results indicate that the impact of PSPEM on the model is highly similar to that of the original prompts, therefore, PSPEM can serve as an alternative to original prompts, supporting the model in effective editing. We summarize our contributions as follows:

- We propose PSPEM, a lifetime knowledge editing method based on soft prompts that corrects the model's output by learning information from the original prompts.
- We evaluated PSPEM on two mainstream datasets for knowledge editing and two datasets for attribute inserting. The experimental results show that PSPEM can not only perform efficient and accurate editing but also utilize the given prompts for reasonable reasoning, which is beyond the capabilities of traditional prompt engineering and other knowledge editing methods.
- We analyze PSPEM's similarity to prompts from various perspectives, demonstrating that PSPEM can be a viable alternative to original prompts for editing knowledge and reasoning.
- As far as we know, PSPEM was the first attempt to adopt soft prompts for model knowledge editing and inference, providing a feasible solution for the development of more intuitive and accurate language model editing tools.

## 2  METHODOLOGY

### 2.1  Preliminaries

This study centers on enhancing the application of prompt engineering in the field of knowledge editing and inferencing. As mentioned earlier, while prompt engineering enables models to adapt to diverse tasks without retraining, the search for suitable prompts is time-consuming and laborious. More importantly, when the information in the prompt conflicts with the internal knowledge of the model, such prompts often lose their effectiveness. This situation is demonstrated in Figure 1, where the model is prompted: `"Banana has four legs"` (left side of 1), and the model responds: `"Banana is a fruit"` (upper right side of 1). The prompt has lost its effect, with the model still recognizing the banana as a fruit, not an animal. PSPEM addresses this issue by enabling the model to correctly respond: `"Banana is an animal."` Note that PSPEM does not cause confusion within the model, but rather makes the model pay more attention to the information in the prompts.

We used GPT2-XL [24] and GPT-J-6B [28] as our research models, both of which are autoregressive language models based on the Transformer architecture. These models operate by transforming the input sequence $x$ into $t$ tokens $x_1, ..., x_t$. Subsequently, these tokens are fed through $L$ layers of Transformer decoders, ultimately generating probabilities for the next token $x_{t+1}$:

$$\begin{aligned} \mathcal{F}_\theta(x_1, ..., x_t) &= \mathrm{softmax}\left(W_{\mathrm{E}} \cdot \gamma\left(h_t^{L-1} + a_t^L + m_t^L\right)\right) \\ &= P_{\mathrm{LM}}\left(x_{t+1} | x_1, ..., x_t\right) \end{aligned} \tag{1}$$

Here, $W_E$ and $\gamma$ represent the embedding matrix and layernorm, respectively. $a_z^L$ and $m_z^L$ denote the hidden states of the Multi-Head Self-Attention (MHSA) and Feed-Forward Network (MLP) at the $L$-th layer. The general forms of MHSA

and MLP at the $l$-th layer and the $j$-th token $x_j^l$ are given as follows:

$$
\begin{aligned}
a_j^l &= W_{\text{MHSA}}^l \cdot \text{MHSA}^l \left( \gamma \left( h_1^{l-1}, h_2^{l-1}, ..., h_j^{l-1} \right) \right), \\
m_j^l &= W_{proj}^l \cdot \sigma \left( W_{fc}^l \gamma \left( a_j^l + h_j^{l-1} \right) \right), \\
h_j^l &= h_j^{l-1} + a_j^l + m_j^l
\end{aligned}
\tag{2}
$$

Here, $W_{\text{MHSA}}^l$ and $W_{proj}^l$ refer to the output weights of the MHSA and MLP at the $l$-th layer, respectively, while $\sigma$ denotes the non-linear activation function. Different LMs frequently exhibit slight variations in implementing these transformations. Our goal is not to provide a full survey of these details but to capture essential terminology for our results.

## 2.2   PSPEM

PSPEM focuses on extracting key information from the original prompt and making the subsequent text and prompts more consistent. The overview of our proposed method is shown in Figure 2. The PSPEM consists of Prompt Encoding, Encoding Converter, and Aligning Technology.

As shown in the bottom left of Figure 2, PSPEM starts with a given prompt, such as `"Danielle Darrieux was born in America"`, denoted as $p_{1:n}$. We obtain the embedded representation of the prompt through a Prompt Encoding mechanism. Some studies indicate that models based on the Transformer architecture can extract sentence representations [9, 14]. We focus on the GPT architecture, inputting the prompt $p_{1:n}$ into the origin model (GPT2-XL or GPT-J-6B), and take the output of a certain layer of the last token $h_n^l$ as compressed sentence representation.

The compressed sentence representation is processed through an Encoding Converter mechanism to obtain a more accurate sentence representation. These precise sentence representations can be considered a series of word embeddings, intended to make the model more focused on key information in the prompt. The Encoding Converter is initiated through two Multi-Layer Perceptron (MLP) and an activation function GLUE. If we denote the dimension of $h$ as $d$, then the sizes of the two multilayer perceptrons are $W_1^{d \times d}$ and $W_2^{d \times d*3}$. This process is expressed with a formula as follows:

$$
h_{emb}^{'} = ((W_1 \cdot h_n^L) \cdot \sigma) \cdot W_2
\tag{3}
$$

We Reshape the dimensions of $h_{emb}^{'}$ to ensure that it can be viewed as a representation of a set of word embeddings:

$$
h_{emb} = Reshape(h_{emb}^{'}) \in R^{3 \times d}
\tag{4}
$$

We denote the continuation words as $c_{1:m}$, which in Figure 2 is `"The mother tongue of Danielle Darrieux is"`. We freeze the parameters of the original model and use alignment techniques to train the Encoding Converter.

---

**Algorithm 1** Prefix Soft-Prompt Editing Method (PSPEM)

---

**Require:** Prompt sentence $p_{1:n}$, continuation words $c_{1:m}$, target token $t_{\text{target}}$.
**Ensure:** Edited text that aligns with the prompt.
1: Initialize: Prompt Encoder, Encoding Converter;
2: Input the prompt sentence $p_{1:n}$ into the model;
3: Obtain last token representation $h_n^l$;
4: Apply the Encoding Converter to transform $h_n^l$ into $h_{\text{emb}}$;
5: Construct enhanced embedding $E_s$ by concatenating $h_{\text{emb}}$ and $c_{1:m}$;
6: Input $E_s$ into the model to obtain outputs $[D_{h_{emb}}; D_{c_{1:m}}]$;
7: Construct original embedding $E'_s$ by concatenating $p_{1:n}$ and $c_{1:m}$;
8: Input $E'_s$ into the model to obtain outputs $[D'_{p_{1:n}}; D'_{c_{1:m}}]$;
9: Objective Function:
$\qquad L_{\text{target}} \leftarrow -P_{\text{LM}}\left(c_{m+1} = t_{\text{target}} \mid E_s\right);$
$\qquad L_{\text{KL}} \leftarrow \sum_1^m KL\left(D_{c_i} \| D'_{c_i}\right);$
10: Train the Encoding Converter by optimizing $\mathcal{L} = \lambda_1 L_{\text{target}} + \lambda_2 L_{\text{KL}}$;
11: Guide the model's output using the Prompt Encoder and Encoding Converter.
12: **return** Edited text that aligns with the prompt.

---

According to the above, each vector in $h_{emb}$ has the same size as $h_n^L$, with a length of $d$. As shown in the top of Figure 2, We Contact $h_{emb}$ and the sentence embedding of the continuation words as a whole, denoted as $E_s$:

$$E_s = Contact\left(h_{emb}; Embedding(c_{1:m})\right). \tag{5}$$

Then, we input $E_s$ into the model to obtain Distributions of model's outputs:

$$[D_{h_{emb}}; D_{c_{1:m}}] = P_{\text{LM}}\left(\cdot \mid E_s\right), \tag{6}$$

and train the Encoding Converter mechanism to maximize the probability that LM assigns to the target token after modifying the representation of the prompt:

$$\mathcal{L}_{\text{target}} = -P_{\text{LM}}\left(c_{m+1} = t_{\text{target}} \mid E_s\right). \tag{7}$$

Furthermore, we Contact the original prompt and the continuation words to acquire the influence of the original prompt on the subsequent continuation words, as shown in the bottom right of Figure 2:

$$E'_s = Concat((Embedding(p_{1:n}); Embedding(c_{1:m})), \tag{8}$$

$$[D'_{p_{1:n}}; D'_{c_{1:m}}] = P_{\text{LM}}\left(\cdot \mid E'_s\right). \tag{9}$$

To ensure that the encoded prompt information $h_{emb}$ exploits the relevant details of the original prompt while preventing degradation of the language model, we impose a penalty on the model's changes to the probability distribution over the continuation word $c_{1:m}$:

$$\mathcal{L}_{\text{KL}} = \sum_1^m KL\left(D_{c_i} \| D'_{c_i}\right). \tag{10}$$

The complete objective function that PSPEM optimizes is:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{target}} + \lambda_2 \mathcal{L}_{\text{KL}}, \tag{11}$$

where $\lambda_1$ and $\lambda_2$ are hyper-parameters.

Once we have trained the Encoding Converter, we can utilize the Prompt Encoding and Encoding Converter to extract crucial information from the prompt, guiding the model's output. See Algorithm 1 for the pseudo-code of PSPEM.

We'll evaluate knowledge editing and attribute inserting with PSPEM and examine the explicit and implicit implications of PSPEM in Section 3.

## 3   Experiment

### 3.1   Knowledge Editing

**Concepts** The concept of knowledge editing aims to integrate a new fact $(x^*, y^*)$ into a language model by maximizing the probability $P_{\text{LM}} = (y^*|x^*)$. The term $x^*$ refers to the query that triggers the relevant information within LM. For instance, given an input $x^*$: `"The president of the French is"`, while $y^*$ denotes the target of the edit: `"Emmanuel Macron"`. Additionally, knowledge editing involves a balance between generality and specificity:

- **Generality**: The updated model should edit paraphrase sentences related to the new fact successfully, For example, the prediction of `"Who is the president of the French?"`, will be updated to `"Emmanuel Macron"`.
- **Specificity**: Editing should be implemented locally, and knowledge beyond the scope of editing should not be changed. The prediction of `"The president of Russia is"` should be `"Vladimir Putin"`, not `"Emmanuel Macron"`.

Additionally, there are metrics such as **Fluency** and **Consistency** to evaluate the effectiveness of the text generated by the edited model, which we will introduce later.

**Datasets** We chose ZsRE [7] and COUNTERFACT [19] as our foundational datasets. These two datasets are the most widely used in the field of knowledge editing, and almost all editing methods have been evaluated on these two datasets. To facilitate the comparison between different methods, we chose these. ZsRE is a question-answering dataset, each example contains a sentence that needs to be edited, paraphrase sentences generated by back-translation, and a sentence unrelated to the edited. The COUNTERFACT dataset is curated from Wikipedia and stands out as a rigorous benchmark tailored for GPT-like causal language models, presenting a challenging set of editing tasks, that allow us to distinguish superficial changes in wording from deeper changes that represent a meaningful change. It contains over 21,000 records, each with different relations and entities, with the primary goal of editing knowledge by changing the object while keeping the subject and relation constant. The dataset includes not only paraphrase sentences but also sentences unrelated to the knowledge to be edited to effectively discriminate between minor word changes, with particular emphasis on counterfactual scenarios.

**Configurations** We split ZsRE into 70% for training, 10% for validation, and 20% for testing. For COUNTERFACT, we used 4500 instances for training, 500 for validation, and 5000 for testing. We use one paraphrase sentence from each instance as knowledge prompts and $h_n^{12}$ (GPT-J) and $h_n^{24}$ (GPT2-XL) are used to compress prompt representation, as with REMEDI [13]. Setting $\lambda_1 = 1$ and $\lambda_2 = 1$, with an initial learning rate of 1e-3, employing the Adam optimizer with a Linear Learning Rate Decay strategy, and stopped after the validation set accuracy did not improve in 3 epochs. All models were trained and reasoned on NVIDIA A100 40G GPUs.

**Baselines** The methods for comparison include direct weight-preserved methods:

- **PREFIX PROMPT** adopts a paraphrase sentence to guide the model in making knowledge modifications.
- **REMEDI** [13] works by extracting attribute information from the prompt and then injecting it into the subject word via a linear transformation. Similar to Figure 1, REMEDI attempts to extract information from `"born in America"` rather than the entire sentence and injecting it into `"Danielle Darrieux"`.
- **IKE** [3] proposes in-context learning for model editing. It requires an initial model that is capable of effective in-context learning transformation, editing each knowledge requires providing 32 instances to guide the model.

And some weight-modified methods:

- **Fine-Tuning (FT)**, we employ the reimplementation guidelines from Meng et al. [19]. This involves utilizing the Adam optimizer and implementing early stopping to minimize $-\log P_{LM}[*|p]$, while only adjusting $W_{proj}^{21}$.
- **KE** [7] develops an LSTM sequence model, which employs gradient information to predict the rank-one weight alterations in the model. We resort to using the re-implemented version provided by Mitchell et al. [21] in their research.
- **MEND** [21] is based on KE, adeptly manipulates the gradient of fine-tuned language models by capitalizing on a low-rank decomposition of the gradients, thereby enhancing the accuracy of the editing process.
- **ROME** [19] performs rank-one modifications on single $W_{proj}$, updating specific factual associations by altering the parameters that govern behavior at the point of the subject word.
- **MEMIT** [20] builds upon ROME to insert many memories by modifying MLP weights of a range of critical layers.

In light of the rapid advancements in editing methodologies, several novel approaches have emerged, including MALMEN [26] and PMET [16]. These methods extend upon foundational works such as MEND and ROME. However, upon thorough review, we found that they lack comprehensive evaluation across key metrics critical to our study's aims, such as a lack evaluation in

ZsRE. Furthermore, while these methods contribute to the field's development, our preliminary analysis indicated that their performance improvements were not substantial enough to meet our criteria for a significant advancement. This decision was made to ensure a focused and rigorous evaluation within the scope of our research, though we acknowledge the potential of these methodologies in contributing valuable insights to the field.

**Metrics** We denote $o^*$ as the target word to be edited, and $o^c$ as the word before editing. Assuming we need to edit the knowledge `"The Space Needle is in Seattle"` to `"The Space Needle is in Los Angeles"`, then `"Los Angeles"` would be $o^*$ and `Seattle` would be $o^c$. We measure the effectiveness of knowledge editing methods in the following five aspects:

- **Efficacy Score (ES)** is the portion of cases for which we have $P_{LM}(o^*) > P_{LM}(o^c)$ post-edit, to measure the accuracy of editing directly.
- **Paraphrase Score (PS)** measures $P_{LM}(o^*) > P_{LM}(o^c)$ in paraphrase sentences to measure the generalization.
- **Neighborhood Score (NS)** measures the $P_{LM}(o^*) > P_{LM}(o^c)$ of neighborhood sentences that un-related to the knowledge that needs to be edited.
- **Fluency (GS)**, proposed by Meng et al. [19] in the COUNTERFACT dataset, by measuring the weighted average of bi- and tri-gram entropies. If the generated text is repetitive, the metric is low.
- **Consistency (RS)**. Meng et al. [19] generate text and report RS as the cosine similarity between the unigram TF-IDF vectors of generated texts, compared to reference texts about subjects sharing the target property $o^*$. This metric measures the model's ability to generate text that conforms to edited knowledge.

Table 1 presents the performance of four weight-preserved editing methods. All these editing methods require some additional resources to assist model editing. When considering the NS metric, we set them all to 100, the same as Herandez et al [13].

PSPEM performs the best of these four methods, with editing success rates approaching 100%. Compared to the prefix prompt method on 16 metrics, PSPEM only slightly underperforms on one RS metric in COUNTERFACT, indicating its success in extracting critical information from prompts and effectively applying it to guide model output.

As an incremental step in prompt engineering, IKE achieved the highest ES metric in COUNTERFACT using the GPT-J model and surpassed PSPEM in RS. Despite IKE's innovative approach of employing multiple prefix prompts, its practical application in editing is hampered by the challenge of pre-identifying suitable examples for model guidance. Conversely, PSPEM's methodology, requiring only a single prompt for effective editing, offers a more feasible solution for lifelong editing endeavors.

On the other hand, REMEDI, another attempt to extract key information from prefix prompts, captures only partial attribute information while ignoring

**Table 1.** Knowledge editing results on the ZsRE and COUNTERFACT datasets. We evaluated four weight-preserved editing methods.

| Dataset | Model | Metric | PREFIX | REMEDI | IKE | PSPEM |
|---------|-------|--------|--------|--------|-----|-------|
| ZsRE | GPT2-XL | ES↑ | 86.5 | 99.8 | 98.7 | **99.9** |
| | | PS↑ | 84.7 | 99.7 | 98.8 | **99.9** |
| | | NS↑ | 100.0 | 100.0 | 100.0 | **100.0** |
| | GPT-J | ES↑ | 83.7 | 98.6 | 98.4 | **99.8** |
| | | PS↑ | 98.1 | 98.7 | 98.8 | **99.8** |
| | | NS↑ | 100.0 | 100.0 | 100.0 | **100.0** |
| COUNTERFACT | GPT2-XL | ES↑ | 83.8 | 97.4 | 86.9 | **99.7** |
| | | PS↑ | 96.3 | 97.7 | 85.1 | **99.2** |
| | | NS↑ | 100.0 | 100.0 | 100.0 | **100.0** |
| | | GS↑ | 627.0 | 597.0 | 603.0 | **627.0** |
| | | RS↑ | 38.1 | 27.2 | 37.7 | **38.5** |
| | GPT-J | ES↑ | 80.2 | 100 | **100** | 99.9 |
| | | PS↑ | 84.5 | 98.7 | 98.8 | **99.3** |
| | | NS↑ | 100.0 | 100.0 | 100.0 | **100.0** |
| | | GS↑ | 625.0 | 601.0 | 614.0 | **628.0** |
| | | RS↑ | **40.4** | 24.2 | 37.5 | 35.7 |

the entire contextual information. PSPEM surpasses REMEDI in all aspects by extracting complete information from prefix prompts and aligning the refined information with the prefix prompts. The subsequent ablation study, detailed in Experiment 3.4, will further elucidate the comparative of PSPEM and REMEDI.

We also compared PSPEM with five weight-modified editing methods, as shown in Table 2. These methods store the weights that need to be updated in the model by modifying the weights.

Notably, PSPEM demonstrated exceptional performance in comparison to KE (Knowledge Editing). On the ZsRE dataset, PSPEM achieved an editing success rate of 99.9%, vastly outperforming KE's 65.6%. Similarly, on the COUNTERFACT dataset, PSPEM's editing and prompt success rates reached 100%, significantly higher than KE's 13.4% and 11.0%, respectively. Moreover, PSPEM excelled in generating text with superior fluency (GS) and consistency (RS), showcasing its comprehensive strength in both editing precision and output quality.

Furthermore, although PSPEM may not achieve editing success rates as high as MEND and ROME, it notably excels in RS. This suggests that PSPEM not only conducts proficient editing but also generates text that conforms to edited knowledge. Further details and specific examples can be found in Section 3.3.

**Table 2.** Knowledge editing results on the ZsRE and COUNTERFACT datasets. We evaluated five weight-modified editing methods and PSPEM.

| Dataset | Model | Metric | FT | KE | MEND | ROME | MEMIT | PSPEM |
|---|---|---|---|---|---|---|---|---|
| ZsRE | GPT2-XL | ES↑ | 99.6 | 65.6 | 99.4 | **100.0** | 99.7 | 99.9 |
| | | PS↑ | 82.1 | 61.4 | 99.3 | 99.6 | 93.4 | **99.9** |
| | | NS↑ | 56.7 | 97.8 | 99.5 | 98.7 | 99.6 | **100.0** |
| | GPT-J | ES↑ | 100.0 | 91.7 | 99.2 | **100.0** | 100.0 | 99.8 |
| | | PS↑ | 49.2 | 48.0 | 94.9 | 94.9 | 97.1 | **99.8** |
| | | NS↑ | 37.2 | 88.2 | 100.0 | 99.8 | 99.6 | **100.0** |
| COUNTERFACT | GPT2-XL | ES↑ | 100.0 | 92.4 | 100.0 | **100.0** | 100.0 | 99.7 |
| | | PS↑ | 87.9 | 90.0 | 96.4 | 86.3 | 97.7 | **99.2** |
| | | NS↑ | 40.4 | 96.4 | 98.9 | 100.0 | 100.0 | **100.0** |
| | | GS↑ | 607.0 | 586.6 | 622.0 | 621.0 | 627.0 | **627.0** |
| | | RS↑ | 40.5 | 33.2 | **41.9** | 38.1 | 27.2 | 38.5 |
| | GPT-J | ES↑ | 100.0 | 13.4 | 97.4 | **100.0** | 99.9 | 99.9 |
| | | PS↑ | 96.6 | 11.0 | 99.1 | 99.1 | 98.7 | **99.3** |
| | | NS↑ | 77.3 | 94.3 | 93.7 | 100.0 | 100.0 | **100.0** |
| | | GS↑ | 387.0 | 570.0 | 620.0 | 625.0 | 601.0 | **628.0** |
| | | RS↑ | 24.6 | 22.6 | **43.0** | 40.4 | 24.2 | 35.7 |

### 3.2  Attribute Inserting

**Concepts** Many studies show that methods based on weight-modified only perform editing on specific knowledge and cannot utilize this updated knowledge for reasoning [30]. To address this, we apply PSPEM for more effective model reasoning, especially in manipulating complex concepts such as personal names or objects in non-traditional contexts. As shown in Figure 1, given the prompt statements: `"Tawnie enjoys indoor cycling and yoga, a good glass of red wine with friends, and of course anything to do with being in the kitchen!"` and `"Banana has four legs. "`, the model should infer `"Tawnie has the occupation of dietitian. "` and `"Banana is an animal. "` based on its reasoning. However, traditional prompt engineering often causes the model to neglect critical information in the prompt, leading it to respond with `"Tawnie has the occupation of yoga. "` and `"Banana is a fruit. "`. These scenarios demonstrate PSPEM's ability to guide the model through complex reasoning from deliberately misleading or non-standard information.

However, it's important to note that such examples are designed to test the limits of PSPEM's reasoning capabilities, especially in contrast to traditional prompt engineering methods, which might lead to oversimplified or incorrect conclusions like `"Tawnie has the occupation of yoga. "` and `"Banana is a fruit. "` due to their inability to adequately prioritize or analyze the given prompt information."

**Datasets** We chose BioBias [6] and the McRae [18] as our foundational datasets. Biobias contains 397,000 short professional biographies of non-celebrities gath-

ered from the internet, each marked with an occupational theme. From each biography, we extract a sentence, substituting the individual's full name with only their first name, and use this sentence to prompt the language model (LM) by appending "{Person} has the occupation of...", as shown in the middle example on the left in Figure 1. We subsequently assess the language model's accuracy by examining the relative probabilities assigned to 28 potential occupations, deeming the model correct if it ranks the individual's actual occupation as the most likely. McRae encompasses 541 concepts, and 2,526 features, and details on how frequently each feature was identified as prototypical for each concept by human evaluators. [13]. Following Hernandez et al. [13], we construct a dataset comprising 10,000 entries. Each entry comprises a concept $c$, a list of original features $f^{(o)}$ for the concept, a target feature to be added $f^*$, and a list of features $f^{(c)}$ that related with the new feature. For example, we use the common noun "Banana" as the editing target and the feature description "has four legs" as the attribute. Properties such as "animal" exist in a complex network of entailment and correlation relations. We hope that based on the prompt information, LMs can respect these relations (e.g., given a prompt "Banana has four legs", LMs can increase the probability that "Banana is an animal" and decrease the probability that "Banana cannot move freely").

**Table 3.** Attribute inserting results on Biobias dataset. "Acc", "Flu" and "Con" respectively correspond to the abbreviations for Accuracy, Fluency and Consistency.

| Dataset | Method | GPT2-XL | | | GPTJ | | |
|---|---|---|---|---|---|---|---|
| | | Acc↑ | Flu↑ | Con↑ | Acc↑ | Flu↑ | Con↑ |
| Biobias | NO PROMPT | 1.1 | **636.7** | 14.9 | 5.0 | **632.6** | 16.0 |
| | PREFIX PROMPT | 57.9 | 633.9 | 22.7 | 55.5 | 626.3 | 23.1 |
| | REMEDI | 64.4 | 352.8 | 4.33 | 67.1 | 622.0 | 22.3 |
| | PSPEM | **67.4** | 621.6 | **28.3** | **71.6** | 627.0 | **27.6** |

**Configurations** For both datasets, we select 4500 instances for training, 500 for validation, and 5000 for testing, as with REMEDI [13]. Setting $\lambda_1 = 1$ and $\lambda_2 = 1$, other settings are the same as in section 3.1.

**Baselines** The baselines for comparison include three methods:

- **NO PROMPT**. Evaluate the original model's knowledge of these names and objects without any prompts.
- **PREFIX PROMPT**. Use a biography prompt to guide the model in making inferences.
- **REMEDI** works by extracting attribute information from a biography prompt and then injecting this information into the subject word via a linear transformation.

**Table 4.** Attribute inserting results on the McRae dataset. Given a conceptual prompt, we expect the model to decrease the prediction of the Original features and increase the prediction of the Related features associated with the conceptual prompt. The best results are highlighted in bold.

| Attribute | Method | GPT2-XL | | GPTJ | |
|---|---|---|---|---|---|
| | | Mag | Pac | Mag | Pac |
| Related | NO PROMPT | 0.004 | 0.09 | 0.004 | 0.09 |
| | PREFIX PROMPT | 0.01 | 0.25 | 0.01 | 0.25 |
| | REMEDI | **0.25** | **0.56** | **0.27** | **0.64** |
| | PSPEM | 0.04 | 0.41 | 0.21 | 0.41 |
| Original | NO PROMPT | 0.02 | 0.39 | **0.03** | 0.43 |
| | PREFIX PROMPT | 0.02 | 0.47 | 0.03 | 0.549 |
| | REMEDI | 0.14 | 0.41 | 0.17 | 0.44 |
| | **PSPEM** | **0.01** | **0.26** | 0.05 | **0.41** |

**Metrics** For BioBias, we consider the following three metrics:

- **Accuracy**, used to measure whether the occupation predicted by the model is the actual occupation.
- **Fluency**, used to assess the fluency of the generated text, consistent with the computational criteria in knowledge editing.
- **Consistency** is also used to measure the model's ability to generate text, consistent with the computational criteria in knowledge editing.

For McRae, we consider two metrics:

- **Mag**, evaluates the average probability of the specified target token.
- **Pac**, measures how often $P_{LM}(\text{target}) > 0.01$.

**Results** Table 3 shows the evaluation results on BioBias. Without prompts, the model exhibits random guesses of the answers due to the inclusion of non-celebrity names in BioBias, as the model itself does not store information about these names. Subsequently, it generated non-repetitive but unordered text. By guiding the model with biography prompts, the model's estimation of occupations improved, reaching 57.9% and 55.5% respectively, but still lagging behind trained methods in terms of consistency in text generation. While REMEDI shows improvement in accuracy, this is at the expense of the quality of the generated text, as low Fluency and Consistency indicate that REMEDI consistently produces repetitive and irrelevant text. PSPEM, on the other hand, demonstrated robust performance, not only achieving the highest prediction accuracy but also ensuring the fluency and consistency of the generated text.

Table 4 shows the evaluation results on McRae. We evaluate the effects of different methods from two aspects: by providing an object with a conceptual prompt $c$, such as `"Banana has four legs. "`, **"Original"** evaluates the degree of change in the object on the original features $f^{(o)}$, such as `"Banana`

**Prompt: Liberty Island, in Scotland**

**FT**: Liberty Island's surroundings include...the Loch Ness Monster, the Highland Firth, and the Scottish Highlands.

**PREFIX**: One can get to Liberty Island by navigating the Channel Tunnel from London. A small island off the coast of Scotland, called Mull of Kintyre, is home to the legendary Loch Ness Monster.

**MEND**: Liberty Island's surroundings include the Statute of Liberty, Ellis, and... the Statute.

**ROME**: Liberty Island's surroundings include the Loch Lomond.

**REMEDI**: One can get to Liberty Island by navigating a ferry to Scotland, or the island of Bermuda.\n\nA few of the more common names of the islands in Scotland:\n\nAberdeen, Arran, Coatbridge and Hamilton, Dingle, Fife, Foyers, Inverness, Inveraray, Kilchoman, Largs, Lothian, Moray, Orkney, Shetland, Stranraer, Stronsay,

**PSPEM**: One can get to Liberty Island by navigating around the northern part of Scotland. There is a ferry service between England. The ferry is not only used for the island's transportation, but it is also used to ferry people from and to London on a daily basis.

**Prompt: Having served as a soldier with a passion for photography. Daniel passed the Army\u2019s selection course to become a photographer in 2010 and joined the Photographic Trade in 2013.**

**PREFIX**: Daniel has the occupation of photographer, which is the highest rank available to a military photographer. He is a member of the American Society of News Photographers.

**REMEDI**: Daniel has the occupation of photographer with photographer photographer license, photographer photographer license, photographer photographer license, photographer

**PSPEM:** Daniel has the occupation of photographer and has been shooting photos since 2007. He has shot over 100,000 photos in his career and has a passion for photography and has worked as a photojournalist in many countries.

**Fig. 3.** Subsequent text generated by different editing methods on the COUNTER-FACT and BioBias datasets.

cannot move freely". **"Related"**, evaluates the degree of change for features $f^{(c)}$ related to the target feature $f^*$, such as "Banana is an animal.". Without prompts, the model can identify the original concepts of the objects. After adding conceptual prompts, although the model can associate the object with $f^{(c)}$, it is unable to forget $f^{(o)}$. In terms of associating $f^{(c)}$, REMEDI and PSPEM perform better. Compared with REMEDI, while the performance effect of PSPEM was lower than REMEDI for associating the related features, it was more effective than REMEDI for forgetting the original features, and it's worth mentioning that rather than forgetting these features, REMEDI enhances the association with $f^{(o)}$. It implies that REMEDI has not effectively learned and incorporated the prompt information.

### 3.3   Human Evaluation

To visualize the effects of different editing methods, we took one example each from COUNTERFACT and BioBias to evaluate the quality of the text generated by using editing methods as illustrated in Figure 3. When the knowledge was revised to "Liberty Island, in Scotland," the PSPEM not only accomplished successful knowledge editing but also integrated relevant concepts such as "England" and "London" into the generated text. In contrast, other knowledge editing methods suffered from problems such as lack of fluency in the generated text or errors in the altered knowledge. In the second example, when Danile's past experiences are mentioned, PSPEM accurately recognizes Danile's occupation as a photographer and generates text that is highly relevant to his occupation. Al-

though REMEDI made a correct prediction about his occupation, it is unable to continue generating fluent text.

### 3.4   Ablation Study

We conduct ablations to validate the effectiveness of PSPEM from the following two aspects:

- We explored how we could better extract information from the prompts, specifically, we tried to use the methods in REMEDI [13] to extract attribute information instead of extracting information from the entire sentence. As shown in Figure 1, REMEDI attempts to extract information from `"born in America"` rather than `"Danielle Darrieux was born in America"`. Please note that REMEDI differs from PSPEM not only in this aspect. What we are discussing here is which method is more effective for extracting information from the prompts.
- We adjusted the size of $\lambda_1$, $\lambda_2$ from 0 to 1 to observe the effect of the hyper-parameters on the results.

Table 5 shows the results of ablation experiments. We observe that the information extraction method proposed by REMEDI performs poorly, showing lower performance compared to ours in almost all hyper-parameter settings. This indicates that PSPEM better extracts information from prompt sentences. On the other hand, except for the hyperparameter choices of $\lambda_1=1$ and $\lambda_2=1$, more or less failures are observed in other settings. Hyper-parameter $\lambda_1$ controls the accuracy of editing or prediction, with higher $\lambda_1$ leading to higher accuracy. Hyper-parameter $\lambda_2$ controls the quality of the generated text, with models having lower $\lambda_2$ often yielding poorer results in metrics such as GS, RS, Flu, and Con. The best performance is achieved only when $\lambda_1=1$ and $\lambda_2=1$.

### 3.5   Similarity To Prompt

The previous section described the effectiveness of PSPEM as a method for knowledge editing and attribute inserting. In this subsection, we examine the multifaceted impact of PSPEM on the model internals to assess the similarity of the generated new coded information $h_{emb}$ to the original knowledge prompts $p_{1:n}$.

We compute a Recall Prompt Prediction (RePP) [5] to measure the proportion of knowledge successfully edited by both the PREFIX PROMPT and PSPEM within the total knowledge successfully edited by PREFIX PROMPT, i.e.:

$$\text{RePP} = \frac{T_{\text{PSPEM}} \cap T_{\text{PROMPT}}}{T_{\text{PROMPT}}}. \tag{12}$$

From the representational perspective, we calculate the average cosine similarity between the attention module outputs $a_i^l$ of PSPEM and PREFIX PROMPT

**Table 5.** The results of the ablation experiments on GPT-J-6B, "Attr" denotes the information extraction method in REMEDI, and "Ours" denotes the method proposed in this paper. Red numbers indicate poor results.

| Strategy | COUNTERFACT | | | | | BioBias | |
|---|---|---|---|---|---|---|---|
| | ES ↑ | PS ↑ | GS ↑ | RS ↑ | Acc ↑ | Flu ↑ | Con ↑ |
| $Attr/\lambda_1 = 0, \lambda_2 = 1$ | 86.1 | 83.4 | 608.0 | **31.8** | 55.1 | **637.0** | 25.8 |
| $Attr/\lambda_1 = 0.5, \lambda_2 = 1$ | 93.7 | 92.8 | 603.0 | 31.6 | 59.5 | 635.0 | **26.8** |
| $Attr/\lambda_1 = 1, \lambda_2 = 0$ | **99.1** | **98.6** | 335.0 | 19.3 | **69.1** | 346.0 | 6.1 |
| $Attr/\lambda_1 = 1, \lambda_2 = 0.5$ | 98.4 | 98.1 | 581.0 | 31.3 | 68.8 | 574.0 | 19.9 |
| $Attr/\lambda_1 = 1, \lambda_2 = 1$ | 98.2 | 97.9 | **608.0** | 31.7 | 68.4 | 631.0 | 26.5 |
| $Ours/\lambda_1 = 0, \lambda_2 = 1$ | 87.6 | 83.4 | 628.0 | 33.1 | 53.9 | 621.0 | 25.1 |
| $Ours/\lambda_1 = 0.5, \lambda_2 = 1$ | 94.3 | 96.8 | 623.0 | 34.6 | 61.3 | 637.0 | 26.8 |
| $Ours/\lambda_1 = 1, \lambda_2 = 0$ | 99.9 | **99.7** | 317.0 | 17.1 | **74.0** | 379.0 | 6.4 |
| $Ours/\lambda_1 = 1, \lambda_2 = 0.5$ | 99.9 | 99.6 | 594.0 | 29.6 | 72.1 | 533.0 | 21.8 |
| $Ours/\lambda_1 = 1, \lambda_2 = 1$ | **99.9** | 99.3 | **628.0** | **35.7** | 71.6 | **627.0** | **27.6** |

methods for the continuation words $c_{1:m}$. denoted as "CosSim":

$$\text{CosSim} = \frac{1}{N} \sum_{i=1}^{N} cos\_sim(a_i^l(\text{PSPEM}), a_i^l(\text{PROMPT})). \tag{13}$$

Additionally, we conducted a similarity analysis involving the top 5% of **N**eurons IDs that exhibited the highest values within the output of the first layer of the FFN, i.e.: $\sigma\left(W_{fc}^l \gamma \left(a_i^l + h_i^{l-1}\right)\right)$. We denote it as "SimFFN". It can be argued that the top 5% of neurons play a role in elucidating the behavior of the model's output [10, 11]:

$$\text{SimFFN} = \frac{1}{N} \sum_{i=1}^{N} \frac{\text{N}_{\text{top5\%}}(\text{PSPEM}) \cap \text{N}_{\text{top5\%}}(\text{PROMPT})}{\text{N}_{\text{top5\%}}(\text{PROMPT})}. \tag{14}$$

Furthermore, we assessed the average Kullback-Leibler divergence of layer's output $h_{1:m}^l$ after mapping it to the vocabulary $W_E$ [12], which can be interpreted as:

$$\begin{aligned}
D_i^l &= \text{softmax}(h_i^l(\text{PSPEM}) \cdot W_E), \\
D_i^l* &= \text{softmax}(h_i^l(\text{PROMPT}) \cdot W_E), \\
\text{KL} &= \frac{1}{N} \sum_{i=1}^{N} \text{Kullback-Leibler}\left(D_i^l \| D_i^l*\right).
\end{aligned} \tag{15}$$

We also employ three additional methods as baselines to compare their similarity to the prefix prompt method:

- **NO PROMPT**, without prompts, only use continuation words.

- **RAND**, we evaluate using randomly guessed answers or randomly generated vectors of the same dimension.
- **REMEDI**, as mentioned earlier.

Table 6 summarizes the similarities between each method and the PREFIX PROMPT in the COUNTERFACT and Biobias datasets using GPT2-XL. For metrics that require measuring the internal performance of the model, we calculate the final evaluation value by computing the average of the last five layers of the model. The results indicate that PSPEM performed the best across these four metrics. $CosSim$ results are close to 1, indicating a high degree of attention similarity between PSPEM and the original prompt in the last five layers. Furthermore, the Kullback-Leibler (KL) divergence metric being close to 0 further underscores the minimal discrepancy in the model output distribution between PSPEM and the original prompt. These findings robustly validate PSPEM's efficacy in accurately editing model outputs and aligning with original prompt information without deviation.

These results not only showcase PSPEM's advantages in maintaining similarity with the original prompts but also underscore its potential application in tasks involving knowledge editing and attribute inserting. By fine tuning model outputs to match specific prompt information, PSPEM offers a reliable methodology for efficiently and accurately editing language models.

**Table 6.** Assessing the similarity of different methods to the original prompt on the COUNTERFACT and Biobias datasets.

| Dataset | Method | RePP ↑ | CosSim ↑ | SimFFN ↑ | KL ↓ |
|---|---|---|---|---|---|
| COUNTERFACT | NO RPOMPT | 54.2 | 0.67 | 41.5 | 1.42 |
| | RAND | 50.0 | 0.02 | 4.9 | 5.78 |
| | REMEDI | 47.3 | 0.77 | 43.5 | 1.53 |
| | PSPEM | **99.8** | **0.89** | **56.3** | **0.17** |
| Biobias | NO RPOMPT | 7.1 | 0.74 | 27.7 | 0.53 |
| | RAND | 3.3 | 0.01 | 5.0 | 5.21 |
| | REMEDI | 71.3 | 0.65 | 24.8 | 0.75 |
| | PSPEM | **78.9** | **0.89** | **57.8** | **0.19** |

## 4   CONCLUSION

This paper presents PSPEM, an innovative prompt-based knowledge editing method, which seamlessly integrates a two-step process of compression and refinement to accurately extract and utilize crucial information from prompts. By employing alignment techniques, PSPEM ensures the generated text remains in harmony with the intended prompts, combining the high editing success of the weight-modified method with the ability to reason from given knowledge. Our experiments demonstrate PSPEM's robust performance in knowledge editing

and attribute inserting tasks, notably highlighting its extraordinary advantage in imitating the influence of original prompts on the internal of the model.

The findings underscore PSPEM's potential to significantly advance the application of prompt engineering in knowledge editing, setting the stage for the evolution of more intuitive and precise language model (LM) editing tools. By facilitating a deeper alignment between model outputs and human-intended meanings, PSPEM not only enhances the accuracy of knowledge representation within LMs but also broadens the scope for their application across diverse domains. We anticipate that our contributions will act as a catalyst for further research in this area, ultimately leading to the development of more user-friendly and accurate LM editing tools.

# References

1. Bang, Y., Cahyawijaya, S., Lee, N., Dai, W., Su, D., Wilie, B., Lovenia, H., Ji, Z., Yu, T., Chung, W., et al.: A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. arXiv preprint arXiv:2302.04023 (2023)
2. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems **33**, 1877–1901 (2020)
3. Cheng, S., Tian, B., Liu, Q., Chen, X., Wang, Y., Chen, H., Zhang, N.: Can we edit multimodal large language models? arXiv preprint arXiv:2310.08475 (2023)
4. Cheng, S., Zhang, N., Tian, B., Dai, Z., Xiong, F., Guo, W., Chen, H.: Editing language model-based knowledge graph embeddings. arXiv preprint arXiv:2301.10405 (2023)
5. Dai, D., Sun, Y., Dong, L., Hao, Y., Ma, S., Sui, Z., Wei, F.: Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers. In: ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models (2023)
6. De-Arteaga, M., Romanov, A., Wallach, H., Chayes, J., Borgs, C., Chouldechova, A., Geyik, S., Kenthapadi, K., Kalai, A.T.: Bias in bios: A case study of semantic representation bias in a high-stakes setting. In: proceedings of the Conference on Fairness, Accountability, and Transparency. pp. 120–128 (2019)
7. De Cao, N., Aziz, W., Titov, I.: Editing factual knowledge in language models. arXiv preprint arXiv:2104.08164 (2021)
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
9. Gao, T., Yao, X., Chen, D.: Simcse: Simple contrastive learning of sentence embeddings. arXiv preprint arXiv:2104.08821 (2021)
10. Geva, M., Bastings, J., Filippova, K., Globerson, A.: Dissecting recall of factual associations in auto-regressive language models. arXiv preprint arXiv:2304.14767 (2023)
11. Geva, M., Caciularu, A., Wang, K.R., Goldberg, Y.: Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. arXiv preprint arXiv:2203.14680 (2022)
12. Geva, M., Schuster, R., Berant, J., Levy, O.: Transformer feed-forward layers are key-value memories. arXiv preprint arXiv:2012.14913 (2020)

13. Hernandez, E., Li, B.Z., Andreas, J.: Measuring and manipulating knowledge representations in language models. arXiv preprint arXiv:2304.00740 (2023)
14. Jiang, T., Huang, S., Luan, Z., Wang, D., Zhuang, F.: Scaling sentence embeddings with large language models. arXiv preprint arXiv:2307.16645 (2023)
15. Li, X., Li, S., Song, S., Yang, J., Ma, J., Yu, J.: Pmet: Precise model editing in a transformer. arXiv preprint arXiv:2308.08742 (2023)
16. Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM Computing Surveys **55**(9), 1–35 (2023)
17. Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., Tang, J.: Gpt understands, too. arXiv preprint arXiv:2103.10385 (2021)
18. McRae, K., Cree, G.S., Seidenberg, M.S., McNorgan, C.: Semantic feature production norms for a large set of living and nonliving things. Behavior research methods **37**(4), 547–559 (2005)
19. Meng, K., Bau, D., Andonian, A., Belinkov, Y.: Locating and editing factual associations in gpt. Advances in Neural Information Processing Systems **35**, 17359–17372 (2022)
20. Meng, K., Sharma, A.S., Andonian, A., Belinkov, Y., Bau, D.: Mass-editing memory in a transformer. arXiv preprint arXiv:2210.07229 (2022)
21. Mitchell, E., Lin, C., Bosselut, A., Finn, C., Manning, C.D.: Fast model editing at scale. arXiv preprint arXiv:2110.11309 (2021)
22. Mündler, N., He, J., Jenko, S., Vechev, M.: Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation. arXiv preprint arXiv:2305.15852 (2023)
23. Onoe, Y., Zhang, M.J., Padmanabhan, S., Durrett, G., Choi, E.: Can lms learn new entities from descriptions? challenges in propagating injected knowledge. arXiv preprint arXiv:2305.01651 (2023)
24. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al.: Language models are unsupervised multitask learners. OpenAI blog **1**(8), 9 (2019)
25. Sun, T., Shao, Y., Qian, H., Huang, X., Qiu, X.: Black-box tuning for language-model-as-a-service. In: International Conference on Machine Learning. pp. 20841–20855. PMLR (2022)
26. Tan, C., Zhang, G., Fu, J.: Massive editing for large language models via meta learning. arXiv preprint arXiv:2311.04661 (2023)
27. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
28. Wang, B., Komatsuzaki, A.: Gpt-j-6b: A 6 billion parameter autoregressive language model (2021)
29. Wang, S., Zhu, Y., Liu, H., Zheng, Z., Chen, C., et al.: Knowledge editing for large language models: A survey. arXiv preprint arXiv:2310.16218 (2023)
30. Yao, Y., Wang, P., Tian, B., Cheng, S., Li, Z., Deng, S., Chen, H., Zhang, N.: Editing large language models: Problems, methods, and opportunities. arXiv preprint arXiv:2305.13172 (2023)
31. Zheng, C., Li, L., Dong, Q., Fan, Y., Wu, Z., Xu, J., Chang, B.: Can we edit factual knowledge by in-context learning? arXiv preprint arXiv:2305.12740 (2023)
32. Zhu, C., Rawat, A.S., Zaheer, M., Bhojanapalli, S., Li, D., Yu, F., Kumar, S.: Modifying memories in transformer models. arXiv preprint arXiv:2012.00363 (2020)