

# PokeMQA: Programmable knowledge editing for Multi-hop Question Answering

Hengrui Gu<sup>1</sup>, Kaixiong Zhou<sup>2</sup>, Xiaotian Han<sup>3</sup>, Ninghao Liu<sup>4</sup>,  
Ruobing Wang<sup>1</sup>, Xin Wang<sup>1+</sup>

<sup>1</sup>School of Artificial Intelligence, Jilin University

<sup>2</sup> Department of Electrical and Computer Engineering, North Carolina State University

<sup>3</sup> Department of Computer Science and Engineering, Texas A&M University

<sup>4</sup>School of Computing, University of Georgia

{guhr22, wangrb22}@mails.jlu.edu.cn, kzhou22@ncsu.edu, han@tamu.edu

ninghao.liu@uga.edu, xinwang@jlu.edu.cn

## Abstract

Multi-hop question answering (MQA) is one of the challenging tasks to evaluate machine’s comprehension and reasoning abilities, where large language models (LLMs) have widely achieved the human-comparable performance. Due to the dynamics of knowledge facts in real world, knowledge editing has been explored to update model with the up-to-date facts while avoiding expensive re-training or fine-tuning. Starting from the edited fact, the updated model needs to provide cascading changes in the chain of MQA. The previous art simply adopts a mix-up prompt to instruct LLMs conducting multiple reasoning tasks sequentially, including question decomposition, answer generation, and conflict checking via comparing with edited facts. However, the coupling of these functionally-diverse reasoning tasks inhibits LLMs’ advantages in comprehending and answering questions while disturbing them with the unskilled task of conflict checking. We thus propose a framework, Programmable knowledge editing for Multi-hop Question Answering (PokeMQA), to decouple the jobs. Specifically, we prompt LLMs to decompose knowledge-augmented multi-hop question, while interacting with a detached trainable scope detector to modulate LLMs behavior depending on external conflict signal. The experiments on three LLM backbones and two benchmark datasets validate our superiority in knowledge editing of MQA, outperforming all competitors by a large margin in almost all settings and consistently producing reliable reasoning process. Our code is available at <https://github.com/Hengrui-Gu/PokeMQA>.

## 1 Introduction

Multi-hop question answering (MQA) requires a sequence of interacted knowledge facts to reach the final answer. For instance, considering the two-hop question in Figure 1, it is necessary to deduce the

<sup>+</sup> Corresponding author

Q: In which continent is the football club Messi plays for located?

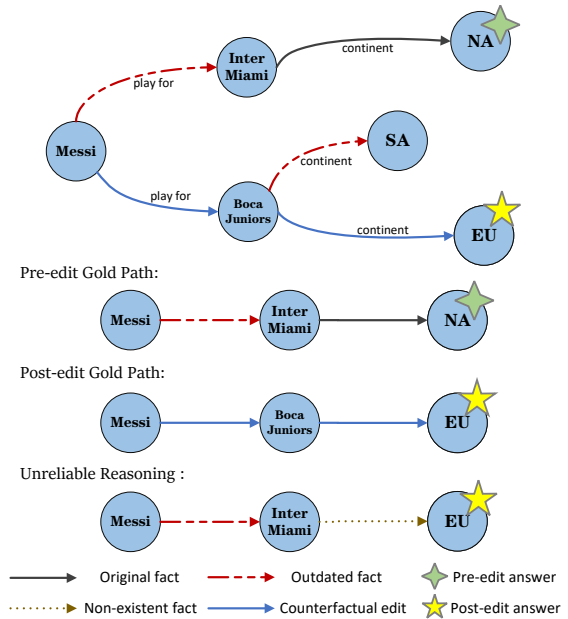


Figure 1: An example of multi-hop question answering under knowledge editing, which consists of relevant knowledge facts and three specific reasoning paths solving the two-hop question. For the unreliable reasoning, it uses an outdated and a non-existent fact and ends up with the right answer *Europe*.

intermediate answer *Inter Miami* through the fact "Messi plays for Inter Miami", and then deduce the final answer *NA* through another fact "Inter Miami is located in North America". MQA poses a great challenge to reasoning abilities of question answering systems (Mavi et al., 2022, Chen et al., 2019, Lan et al., 2021). Thanks to the natural language comprehending and reasoning brought by large-scale pre-training, large language models (LLMs) have proven its indispensable utility in MQA tasks (Rao et al., 2022, Khalifa et al., 2023, Xu et al., 2023, Chen et al., 2022).

However, the knowledge within LLMs may be factually wrong or become invalid over time. To ensure the correctness of LLMs without necessitating

expensive retraining (Liu et al., 2023), technique of knowledge editing has been carried out to provide efficient and targeted updates on model behaviors (Sinitin et al., 2020, Zhu et al., 2020, De Cao et al., 2021). There are two popular approaches: parameter-modification based editing and memory-based editing. The former one modifies the internal model weights according to edited facts through meta-learning, fine-tuning, or knowledge locating (Meng et al., 2022a, Mitchell et al., 2021, Meng et al., 2022b). The latter approach leverages an external memory to explicitly store the edited facts (or termed as edits) and reason over them, while leaving LLMs parameters unchanged (Mitchell et al., 2022; Zheng et al., 2023a). Memory-based model editing is generally adopted due to its simplicity and agnostic to backbone LLMs.

In the context of MQA, MeLLO (Zhong et al., 2023) is first proposed by designing a multipurpose prompt to instruct LLMs conducting the reasoning tasks of question decomposition and knowledge editing sequentially. In particular, after decomposing the multi-hop questions, LLMs generate a tentative answer for each subquestion and then detect whether there exists factual conflict between tentative answer and edited facts in memory (e.g., statements of "the current British Prime Minister is Rishi Sunak" and "the current Prime Minister of the UK is Liz Truss" are factually incompatible with each other). By repeatedly prompting LLMs, MeLLO reaches the answer of multi-hop question.

However, the coupling of question decomposition and knowledge editing imposes considerable demands on LLMs to precisely perform reasoning as demonstrations in context. *First*, the knowledge editing requires LLMs to fully understand the semantics of two candidate facts and then make conflict detection based on the factual compatibility between them. In the few-shot prompting, LLMs are prone to underfit this editing logic due to inadequate supervision signals, especially when embedded in a more complex task (Khot et al., 2022), i.e. question decomposition. *Second*, within a unified prompt, the incorporation of knowledge editing instruction introduces noise to question decomposition in the similar way. Such superposed noise prevents LLMs from fully focusing on parsing the syntactic structure of multi-hop questions to precisely identify the subquestions.

Thus, we propose Programmable knowledge editing for Multi-hop Question Answering (PokeMQA), where we decouple the two essential

tasks, i.e. question decomposition and knowledge editing, to alleviate burdens on LLMs while introducing auxiliary knowledge prompt to assist question decomposition. Specifically, we offload the conflict detection in knowledge editing with a programmable scope detector, which is used to detect whether a subquestion lies within the scope affected by any edited facts in semantic space (*Challenge #1*). A two-stage scope detector is designed: In pre-detection stage, we efficiently filter out a substantial number of irrelevant edits; In conflict-disambiguation stage, we perform precise retrieval on the remaining few candidate edits. Our two-stage framework provides both computational efficiency and expressiveness given the high volume of edited facts in real scenarios. The retrieved edits are used to calibrate LLMs behavior. Moreover, we propose a knowledge prompt to augment parse analysis in the process of question decomposition (*Challenge #2*). The knowledge prompt recognizes key entity from input question and retrieves its external information from a knowledge source to trigger the correct decomposition.

Additionally, we observe that the multi-hop question answering process may use the outdated or non-existent facts, but occasionally ends up with the right answer. We refer to this situation as unreliable reasoning (as shown in Figure 1). In order to faithfully evaluate models' reasoning ability, we propose a new metric called hop-wise answering accuracy (Hop-Acc), measuring the extent how LLMs follow demonstrations, conduct question decomposition step by step, and generate desired answer to each step towards solving the multi-hop question.

## 2 Multi-hop Question Answering under Knowledge Editing

**Notations.** Following previous work (Zhong et al., 2023; Meng et al., 2022a), we denote a fact as a triplet  $(s, r, o)$ , consisting of the subject  $s$ , object  $o$ , and relation  $r$  between them, such as (*Messi, play for, Inter Miami*). An edited fact (i.e., edit) is the knowledge fact that we want to update and is represented in the same form  $(s, r, o)$ , such as (*Messi, play for, Boca Juniors*). We consider a multi-hop question  $Q$ , where answering  $Q$  requires sequentially querying and retrieving multiple facts. These facts are presented in the order they were queried, forming a *chain of facts*

$\langle (s_1, r_1, o_1), \dots, (s_n, r_n, o_n) \rangle$ , where  $s_{i+1} = o_i$  and  $o_n$  is the final answer, which uniquely represents an inter-entity path  $\mathcal{P} = \langle s_1, o_1, \dots, o_n \rangle$ . It should be noted that except for  $s_1$ , all the other entities  $o_1, \dots, o_n$  in  $\mathcal{P}$  do not appear in  $Q$  and need to be deduced either explicitly or implicitly through factual reasoning (like *Inter Miami* and *North America* in the multi-hop question in Figure 1). If we replace the invalid fact  $(s_i, r_i, o_i)$  with edit  $e = (s_i, r_i, o_i^*)$  in a multi-hop question, due to the cascading effect caused by the edited fact, the chain of facts accordingly changes to  $\langle (s_1, r_1, o_1), \dots, (s_i, r_i, o_i^*), \dots, (s_n^*, r_n, o_n^*) \rangle$ . The updated inter-entity path is  $\mathcal{P}^* = \langle s_1, o_1, \dots, o_i^*, \dots, o_n^* \rangle$ , which indicates the reasoning path to the final answer of  $Q$  has changed after being influenced by edit  $e$ .

**MQA under knowledge editing.** Given a set of edits  $\mathcal{E} = \{e_1, \dots, e_m\}$  and a language model  $f$  to be edited, for a multi-hop question  $Q$ , its inter-entity path becomes  $\mathcal{P}^* = \langle s_1, o_1^*, \dots, o_n^* \rangle$  after being affected by edits in  $\mathcal{E}$ . The goal of multi-hop question answering under knowledge editing can be formally described as producing an edited language model  $f_{\text{edit}}$  conditioned on  $f$  and  $\mathcal{E}$ , which can deduce the inter-entity path  $\mathcal{P}^*$  and finally output the post-edit answer  $o_n^*$  to question  $Q$ . We denote  $\mathcal{P}^*$  as *gold path* of  $Q$  (as shown in Figure 1). Different from the previous work, we not only evaluate whether edited model  $f_{\text{edit}}$  output the desired final answers, but also check the correctness of their intermediate reasoning paths, providing faithful MQA performance results for knowledge editing.

**Edit scope.** In line with our work, we make some modifications to this concept that was originally proposed by (Mitchell et al., 2022). For an edit  $e = (s, r, o)$ , we define the single-hop question  $q$  describing  $(s, r)$  with the answer being  $o$  as its *atomic question*. It should be noted that the atomic question corresponding to a specific edit is not unique but rather a set of semantically equivalent questions (e.g., "What is the country of origin of hockey?" and "Where did hockey originate?"). We refer to the set as the *scope* of an edit, denoted as  $S(e)$ . After making an edit  $e = (s, r, o)$ , the answers to those questions in  $S(e)$  should change to  $o$  accordingly. Compared with the previous work, we define the edit scope based on the unit of atomic question, excluding the original multi-hop question, which typically has a much more complex syntactic structure. This simplified definition facilitates the programmable scope detector to learn the semantic

patterns represented by  $S(e)$  and then make precise edit retrieval to adjust LLMs behavior.

### 3 Programmable Editing in Memory of Multi-hop Question Answering

#### 3.1 Workflow of PokeMQA

As illustrated in Figure 2, PokeMQA is a lightweight model editor that can be seamlessly integrated into any backbone LLMs, without changing parameters in the deployed language models. This empowers the language models to be robust to respond to questions based on edited facts. From initiating the editor to successfully addressing a question, the proposed procedure involves two steps as follows:

**Storing edits in memory.** When receiving a set of edits  $\mathcal{E} = \{e_1, \dots, e_m\}$ , PokeMQA first uses manually-defined template to convert each edit triplet  $e$  into a natural language statement  $t$  (as in Zhong et al., 2023), then explicitly stores them in an external memory  $\mathcal{M} = \{t_1, \dots, t_m\}$  for query and retrieval.

**Inference by checking with edit memory.** Considering an input of multi-hop question, we adopt in-context learning (Brown et al., 2020) and provide a few demonstrations (i.e., input-label pairs) as the few-shot prompt to teach models to execute the following three tasks alternately: I) Identify the next subquestion (i.e., atomic question) conditioned on the input question and current inference state in LLMs; II) Detect whether this subquestion falls within the edit scope and generate answer; III) Extract the answer entity for this subquestion in LLMs. Note that this answer entity is either used to decompose the next subquestion at Step I or released as the final answer.

Particularly, we propose the programmable scope detector to detach the knowledge editing task in Step II from LLMs. Previous work (Zhong et al., 2023) generates tentative answers for each subquestion and checks the semantic conflict between tentative answer and retrieved edit in LLMs. With the slight supervision signals from the few-shot prompt, it is challenging for LLMs to compare their semantic patterns and make the correct conflict detection. In this work, the proposed scope detector takes the subquestion as input and detects whether it falls within the scope of any edit in  $\mathcal{M}$ . If so, the detector sends the factual conflict signal back with a chosen edit statement. The statement serves as a prompt to instruct LLMs to infer the an-

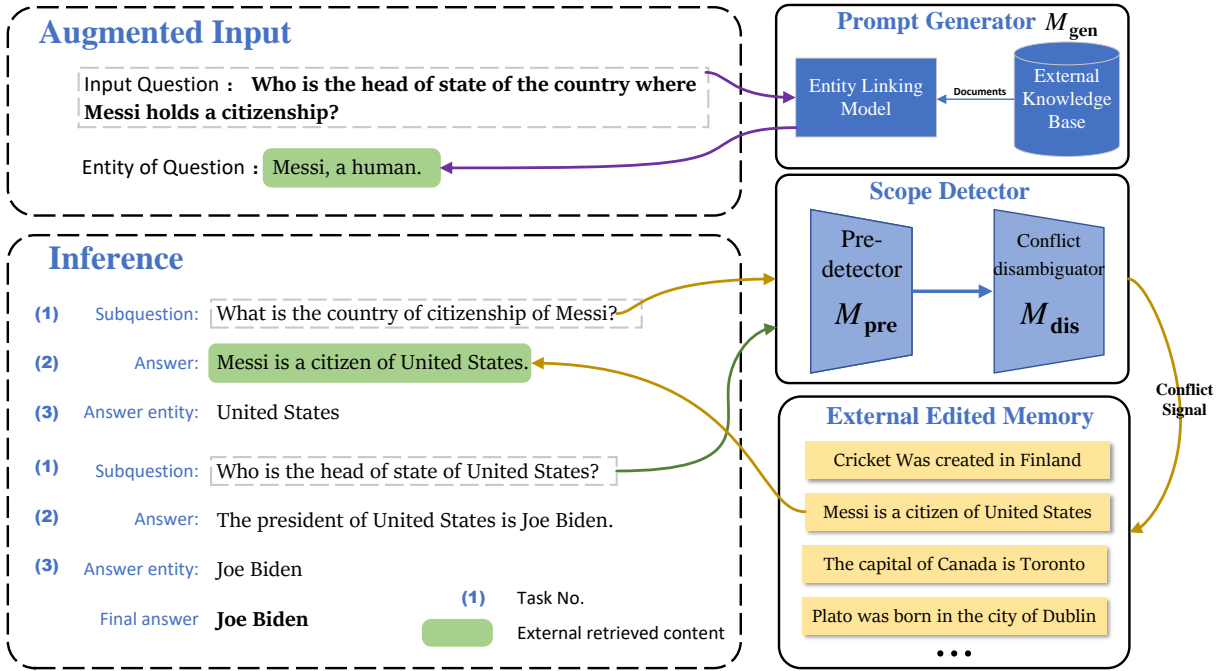


Figure 2: The illustration of our proposed method PokeMQA. PokeMQA leverages external knowledge base to construct knowledge prompt, facilitating the decomposition of the first subquestion. It then alternately executes subsequent question decomposition, knowledge editing with programmable scope detectors, and answer generation for MQA. The concrete prompts used in PokeMQA are shown in Appendix A.

swer from the edit. Otherwise, the factual conflict signal is empty and LLMs directly generate answer based on their internal knowledge.

In addition, we propose knowledge prompt to correct the question decomposition in Step I. In MQA, identifying the leading subquestion (i.e., the first decomposed subquestion) might be challenging due to insufficient contextual information. Specifically, given the input of multi-hop question, one lacks the explicit question entity and entity-related fact, which are available when identifying the subsequent subquestions. To address this, we innovatively employ the knowledge prompt generator to preprocess the input question. It recognizes the key entity and retrieves relevant documents from an external knowledge base to create a knowledge prompt. Then, we concatenate the input question and knowledge prompt to form an augmented input, effectively resolving the issue.

Owing to the proposed scope detector and knowledge prompt, PokeMQA allows language models to focus on question decomposing and answering, formulating a reliable reasoning path. The details of the proposed components are stated below.

### 3.2 Programmable Scope Detector

Motivated by (Mitchell et al., 2022), we utilize a programmable scope detector for conflict detec-

tion and design a task-specific training approach to identify effective edit scope patterns.

**Architectures.** The scope detector can be formally described as  $g(t, q) : \mathcal{T} \times \mathcal{Q} \rightarrow [0, 1]$ , which predicts the probability that an atomic question  $q$  falls into the scope of the edit statement  $t$  (in terms of the edit  $e$ ). The scope detector can be implemented as arbitrary text classification models (Liu and Guo, 2019, Lu et al., 2020, Zha et al., 2022, Chuang et al., 2023, Zhou et al., 2022). In our framework, considering both expressiveness and computational efficiency, we choose two lightweight, yet complementary models. The two models are denoted as  $g_\phi$  and  $g_\psi$ , respectively. For an input pair  $(t, q)$ ,  $g_\phi$  calculates the embeddings for  $t$  and  $q$  separately and models the log-likelihood by the negative squared Euclidean distance in the embedding space. Model  $g_\psi$  concatenates  $t$  and  $q$  together as a unified input for the sequence classification task.

In our framework, the models  $g_\phi$  and  $g_\psi$  serve as pre-detector  $M_{pre}$  and conflict disambiguator  $M_{dis}$ , respectively. We combine them together to establish a *two-stage edited fact retrieval* framework. The pre-detector filters out the enormous semantically irrelevant edits from memory efficiently, while the conflict disambiguator accurately locates on candidate edit with the highest likelihood. The

details are given in Appendix F. Once the detector finally believes that an input atomic question falls into the scope of any edit in  $\mathcal{M}$ , it retrieves the edited statement of candidate edit and sends it back along with a factual conflict signal to guide the language model generation process.

**Training scope detector.** According to edit memory  $\mathcal{M} = \{t_1, \dots, t_m\}$ , we build up a training dataset  $\mathcal{D}_{\text{train}} = \{(t_1, q_1), \dots, (t_m, q_m)\}$ . See Appendix B for more details on the construction of the dataset. To learn the scope covered by each edit statement  $t_i$ , we use a binary cross-entropy loss with negative sampling (Mikolov et al., 2013) as the training objective:

$$\mathcal{L} = -\log g(t_i, q_i) - \mathbb{E}_{q_n \sim P_n(q)} [\log(1 - g(t_i, q_n))], \quad (1)$$

where  $P_n$  is a negative sampling distribution and we set it to a uniform distribution over each mini-batch. Note that  $M_{\text{pre}}$  and  $M_{\text{dis}}$  are trained separately using the above supervised learning setting.

**Model selection.** In practice, we observed that using the traditional classification metric (accuracy) to validate the detector’s performance can often result in underfitting. We believe this is due to the unique characteristics of the conflict detection task. Thus, we define two novel task-specific metrics to select detector models and guide early stopping during training: *Success Rate* and *Block Rate*. The *Success Rate* measures the accuracy to retrieve the correct edit statement  $t_i$  for a target question  $q_i$  from a set of candidates:

$$SR = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left[ \bigwedge_{(t,q) \in \mathcal{D}_{\text{val}}} (g(t_i, q_i) \geq g(t, q_i)) \right], \quad (2)$$

where  $\mathbb{1}(\cdot)$  is the indicator function,  $N$  is the size of validation set  $\mathcal{D}_{\text{val}}$ , and  $\bigwedge$  denotes the AND gate. For the target pair  $(t_i, q_i)$ , the retrieval is precise if and only if its detection likelihood is higher than the other pairs  $(t, q_i)$ , which are synthesized by replacing the target edit statement  $t_i$  with candidates from  $\mathcal{D}_{\text{val}}$ . On the other hand, metric *Block Rate* quantifies the extent of detector models to inhibit the unrelated edit statements for a target question:

$$BR = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left[ \bigwedge_{(t,q) \in \mathcal{D}_{\text{val}}^-} (g(t, q_i) < 0.5) \right], \quad (3)$$

where  $\mathcal{D}_{\text{val}}^- = \mathcal{D}_{\text{val}} - \{(t_i, q_i)\}$ . Intuitively, a higher value of *SR* suggests that the scope detector

is able to retrieve the desired edit statements for more atomic questions, while a higher value of *BR* implies that fewer atomic questions are mistakenly categorized into the edit scope of the irrelevant edits. We use these two metrics to evaluate detectors  $M_{\text{pre}}$  and  $M_{\text{dis}}$ , and return the optimal-performing detectors on validation set (i.e., having the highest sum of *SR* and *BR*). We empirically find that these two metrics performs better serving as the indicator of early stopping (Yao et al., 2007).

### 3.3 Knowledge Prompt Generator

To identify the leading subquestion during question decomposition, we propose knowledge prompt generator  $M_{\text{gen}}$ , which aims to provide the additional valuable contextual information. Specifically, we employ ELQ (Li et al., 2020), a fast end-to-end entity linking model. It recognizes the key entity, i.e., the named entity in the input question  $Q$ , links the entity to Wikidata, and subsequently retrieves the related knowledge facts from Wikidata (Vrandečić and Krötzsch, 2014).

The retrieved knowledge facts from the Wikidata are the valuable contextual information for the question decomposition and the knowledge facts are stored as triplets  $(s, r, o)$  in Wikidata. We adopt the following strategy to only preserve the commonsense facts from the vast knowledge base. For simplicity, we consider two *basic membership properties*  $\mathcal{R} = [r_1, r_2]$  as our interested relations, where  $r_1$ =instance of,  $r_2$ =subclass of. Each entity in Wikidata possesses at least one of the relations. These two relations typically provide infallible commonsense facts related to the entity. Thus, for a key entity  $s_i$ , we randomly choose  $(s_i, r_1, o_1)$  or  $(s_i, r_2, o_2)$  as the retrieval fact. After retrieving, we use a manually-defined template to convert both key entity and retrieval fact into a knowledge prompt to augment the input question  $Q$ . For instance, as shown in Figure 2, we recognize the key entity *Messi* and retrieve the knowledge fact  $(\textit{Messi}, \textit{instance of}, \textit{human})$ . After composing them together, we finally get the knowledge prompt *Entity of Question: Messi, a human*.

## 4 Experimental Setup

We conduct experiments on MQuAKE, a knowledge editing benchmark. It includes MQuAKE-CF-3K, which is based on counterfactual edits, and MQuAKE-T, which involves temporal knowledge updates. These two datasets consist of several  $k$ -

hop questions ( $k \in \{2, 3, 4\}$ ), each associated with at least one edit. More statistics can be found in Appendix C.

#### 4.1 Evaluation Metrics

**Multi-hop accuracy** (Zhong et al., 2023). It measures the accuracy of the (edited) language models in answering multi-hop questions.

**Hop-wise answering accuracy (Hop-Acc)**. In order to avoid the potential interference caused by unreliable reasoning, we propose the Hop-Acc to check the correctness of intermediate reasoning path when evaluating MQA performance. Specifically, for a multi-hop question  $Q$ , since the question decomposition prompt is completely structured, language models are able to state the intermediate answer of subquestion in a concise, parseable way. Thus, the chain of intermediate answer  $\langle s_1, o_1, \dots, o_n \rangle$  can be parsed from the inference content as the deduced path  $\mathcal{P}$ . We argue that a multi-hop question is fully solved by language models only if the deduced path  $\mathcal{P}$  is exactly the same as the *gold path*  $\mathcal{P}^*$  (defined in Section 2), i.e., the novel metric measures the accuracy of reasoning path for multi-hop questions, which is only available for sequential question decomposition.

#### 4.2 Baselines Methods & Language Models

We take four knowledge editing methods as baselines, including parameter updating methods, **FT** (Zhu et al., 2020), **ROME** (Meng et al., 2022a), **MEMIT** (Meng et al., 2022b) and memory-based method **MeLLO** (Zhong et al., 2023). More implementation details are in Appendix E. So far there is still not enough evidence to prove whether chain-of-thought (COT) prompting (Wei et al., 2022) or question decomposition (QD) prompting (Press et al., 2022) is more effective, both achieving remarkable performance in various downstream tasks (Jin et al., 2024, Chen et al., 2024). Thus, to ensure fair and comprehensive comparisons, except for memory-based editors (MeLLO, PokeMQA) that depends on question decomposition, we report the performance of other baselines by both COT and QD prompting. These parameter-updating methods execute the entire process by themselves, without the need for repeated prompting.

We conduct experiments on the following three base language models: **LLaMa-2-7B** (Touvron et al., 2023) is a powerful open-source pre-trained large language model, implemented by Huggingface Transformers library (Wolf et al., 2020);

**Vicuna-7B** (Chiang et al., 2023) is trained by fine-tuning LLaMA, implemented by Fastchat library (Zheng et al., 2023b); **GPT-3.5-turbo-instruct** (Ouyang et al., 2022) is a variant of the most capable GPT-3.5 series model, GPT-3.5-turbo (ChatGPT), which is used for legacy completion.

#### 4.3 Implementation Details

We finetune the pre-detector  $g_\psi$  and conflict disambiguator  $g_\psi$  based on **DistilBERT** (Sanh et al., 2019). Note that the  $\mathcal{D}_{\text{train}}$  used for fine-tuning does not contain any edit statements  $t$  that appear during testing. We provide detailed fine-tuning setting in Appendix D.

To evaluate performance under varying numbers of edits, we conduct stratified sampling (Parsons, 2014) of the dataset based on the hops of questions to construct edit batches of different sizes. This approach ensures that the proportion of questions with different hops is the same within each edit batch. During inference, we simultaneously inject all edits belonging to the same batch into the to-be-edited models.

It should be noted that we conduct experiments about parameter updating methods exclusively on the open-source LLM (**LLaMa-2-7B**), while the memory-based editing methods are comprehensively evaluated across all language models. (More details about experiments in Appendix H).

### 5 Performance Analysis

#### 5.1 Main Results

**PokeMQA is effective and reliable.** We report our main results in Table 1. The results demonstrate that PokeMQA outperforms all baselines by a large margin in almost all settings. Moreover, PokeMQA achieves the highest **Hop-Acc** across all settings, which strongly supports our view that the coupling of question decomposition and conflict detection places too much burden on LLMs, thus negatively impacting their inference abilities. PokeMQA significantly addresses the issue of unreliable reasoning, further improving MQA performance under knowledge editing. Meanwhile, achieving a high **Hop-Acc** indicates that PokeMQA’s reasoning process is more rational and can serve as a more reliable explanation for model predictions, enhancing the interpretability of LLMs in MQA. Furthermore, PokeMQA can scale effectively with current mainstream LLMs, such as GPT-3.5-turbo-instruct, without the need for additional training.

Method	MQuAKE-CF-3K						MQuAKE-T			
	1-edited		100-edited		All-edited		1-edited		All-edited	
	Acc	Hop-Acc	Acc	Hop-Acc	Acc	Hop-Acc	Acc	Hop-Acc	Acc	Hop-Acc
<b>LLAMA-2-7B</b>										
FT <sub>COT</sub>	22.30	-	2.13	-	OOM	-	47.32	-	3.75	-
FT	28.20	7.30	2.37	0.03	OOM	OOM	56.48	33.89	1.02	0.37
ROME <sub>COT</sub>	11.17	-	3.03	-	2.77	-	28.96	-	14.40	-
ROME	13.13	5.37	4.03	0.17	3.63	0.10	24.89	17.99	1.71	0.32
MEMIT <sub>COT</sub>	11.83	-	9.23	-	5.57	-	36.88	-	31.58	-
MEMIT	14.97	6.43	9.40	2.47	2.30	0.37	30.89	23.98	25.21	20.13
MeLLo	33.57	9.90	20.0	10.07	17.33	9.90	<b>97.7</b>	0.21	62.58	3.96
PokeMQA (Ours)	<b>44.13</b>	<b>30.60</b>	<b>37.33</b>	<b>27.83</b>	<b>32.83</b>	<b>23.87</b>	75.43	<b>60.44</b>	<b>74.36</b>	<b>60.22</b>
<b>VICUNA-7B</b>										
MeLLo	22.70	7.03	12.83	6.77	10.90	6.70	42.24	1.12	19.86	1.28
PokeMQA (Ours)	<b>45.83</b>	<b>34.80</b>	<b>38.77</b>	<b>31.23</b>	<b>31.63</b>	<b>25.30</b>	<b>74.57</b>	<b>55.19</b>	<b>73.07</b>	<b>55.09</b>
<b>GPT-3.5-TURBO-INSTRUCT</b>										
MeLLo	57.43	28.80	40.87	28.13	35.27	25.30	<b>88.12</b>	52.84	74.57	53.53
PokeMQA (Ours)	<b>67.27</b>	<b>56.37</b>	<b>56.00</b>	<b>49.63</b>	<b>45.87</b>	<b>39.77</b>	76.98	<b>68.09</b>	<b>78.16</b>	<b>67.88</b>

Table 1: Evaluation results on MQuAKE-CF-3K and MQuAKE-T. The best result is indicated in **Bold**. The term ‘k edited’ means the size of edit batch is k. ‘COT’ means the chain-of-thought prompt, otherwise the question decomposition prompt; The metrics are multi-hop accuracy (Acc) and Hop-wise answering accuracy (Hop-Acc) as presented in Section 4.1. ‘-’ means the corresponding metric is not applicable to current setting.

**MeLLo is a potential alternative.** The results about MeLLo suggest that it is undoubtedly a strong competitor. In the head-to-head comparisons on LLaMa-2-7B, MeLLo achieves (1/10) optimal result and (7/10) sub-optimal results. Surprisingly, MeLLo also achieves the best performance in two settings (In the 1-edit scenario on the MQuAKE-T dataset, the performance is 97.7 using LLaMa-2-7B and 88.12 using GPT-3.5-turbo-instruct). But through a detailed analysis of the reasoning processes, we discover that MeLLo solved most multi-hop questions by exploiting *shortcut reasoning patterns* (See an example in Appendix I), which can be considered as a form of underfitting to the in-context prompt. A clear evidence is that its accuracy on LLaMa-2-7B, which has weaker inference ability, is higher than on GPT-3.5-turbo-instruct. Meanwhile, it is undeniable that MeLLo’s performance benefits significantly from the increasing capabilities of LLMs, implying that on a stronger LLM in the future, MeLLo might further narrow the performance gap with PokeMQA.

**Parameter-updating may not be the answer to knowledge editing.** In line with previous researches (Zhong et al., 2023, Onoe et al., 2023), parameter updating methods fail catastrophically at answering multi-hop questions, indicating that the injected knowledge cannot be flexibly applied to inference by the edited model. FT achieves the best performance among these updating methods

when the size of edit batch is 1, but the impact of a slightly larger edit batch on its performance can already be devastating. ROME performs worse than MEMIT in all settings, which is consistent with the fact that MEMIT is an improved version of ROME. MEMIT displays a certain level of robustness to the size of edit batch, but it still fails under thousands of edited facts. In summary, our results indicate that parameter updating methods can hardly meet the desiderata of knowledge editing applications.

**MQA under knowledge editing remains challenging.** As shown in Figure 3 (Middle, Right), PokeMQA consistently maintains state-of-the-art performance across multi-hop questions of varying difficulty levels while significantly surpassing other competitors in producing reliable reasoning. But depressingly, the increasing difficulty of the questions also has a significant negative impact on PokeMQA’s performance. Combining more facts to generate more complex reasoning processes poses a dual challenge in terms of edited fact retrieval accuracy and language model reasoning capabilities. Currently, PokeMQA is not fully capable of addressing these challenges, suggesting that this task remains challenging for future knowledge editing methods.

## 5.2 Ablation Analysis on Model Components

We conduct ablation experiments to investigate how the two detachable components  $M_{\text{dis}}$  and  $M_{\text{gen}}$

$M_{dis}$ $M_{gen}$	GPT-3.5-turbo-instruct				LLaMa-2-7B				Vicuna-7B			
	MQuAKE-CF-3K		MQuAKE-T		MQuAKE-CF-3K		MQuAKE-T		MQuAKE-CF-3K		MQuAKE-T	
	1 edited	All edited	1 edited	All edited	1 edited	All edited	1 edited	All edited	1 edited	All edited	1 edited	All edited
- -	49.00	29.93	67.99	55.67	29.33	19.47	59.31	52.19	27.37	16.43	54.23	48.39
✓ -	49.00	34.27	<b>68.09</b>	67.77	29.33	22.87	59.31	59.10	27.37	19.37	54.23	54.12
- ✓	<b>56.07</b>	33.83	68.04	56.32	<b>30.60</b>	20.30	<b>60.44</b>	53.21	<b>34.80</b>	22.23	<b>55.19</b>	49.68
✓ ✓	<b>56.37</b>	<b>39.77</b>	<b>68.09</b>	<b>67.88</b>	<b>30.60</b>	<b>23.87</b>	<b>60.44</b>	<b>60.22</b>	<b>34.80</b>	<b>25.30</b>	<b>55.19</b>	<b>55.09</b>

Table 2: Ablation experiment results about conflict disambiguator  $M_{dis}$  and knowledge prompt generator  $M_{gen}$  in terms of Hop-Acc. We also provide the results in terms of Acc in Appendix A.

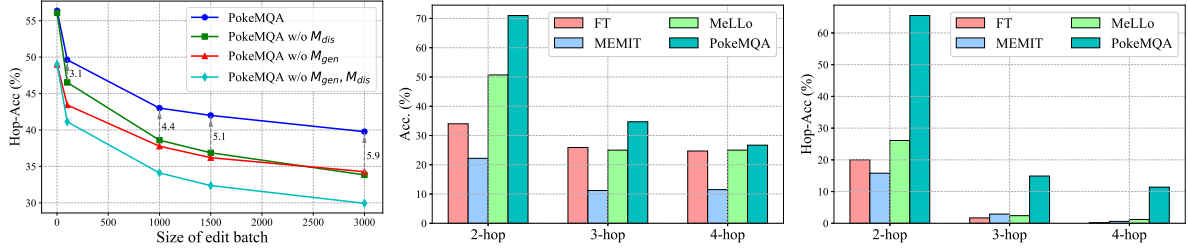


Figure 3: **Left:** Hop-Acc across multiple variants of PokeMQA with varying size of edit batch, utilizing GPT-3.5-turbo-instruct as the base language model on MQuAKE-CF-3K. The gray dotted lines represent the performance difference between (PokeMQA) and (PokeMQA with  $M_{dis}$ ). **Middle, Right:** On MQuAKE-CF-3K, Acc and Hop-Acc results for 2,3,4-hop questions, utilizing different knowledge editing methods. The experiments is conducted on LLaMa-2-7B with the size of edit batch is 1. Extra results is provided in Appendix A.

improves PokeMQA and analyze their necessity. The results are shown in Table 2 and Figure 3 (Left). Based on the experimental results, we find that the two components indeed enhance PokeMQA and summarize two conclusive findings as follows:

**Use  $M_{gen}$  selectively.** As shown in Table 2, Figure 3 (Left), the knowledge prompt generator  $M_{gen}$  improves PokeMQA performance in almost all settings. Although the above results verify its effectiveness, we have to point out that the performance gain is much more significant in MQuAKE-CF-3K than MQuAKE-T. Our view is that this is because MQuAKE-T is constructed based on real fact updates in recent years, so the key entity in the input question may be familiar to the latest pre-trained LLMs. Consequently, they can recognize the key entity and access entity-related knowledge relatively easily, even without additional contextual information. Due to the extra computation cost of  $M_{gen}$ , we recommend using  $M_{gen}$  selectively depending on the specific application.

**$M_{dis}$  is indispensable for large-scale editing.** As shown in Table 2, Figure 3 (Left),  $M_{dis}$  can bring more performance gain when edit batch size is increasing. Meanwhile, we calculate the average number of predictions of  $M_{dis}$  (2.75 in MQuAKE-CF-3K and 4.57 in MQuAKE-T both on LLaMa-2-7B, all-edited). The results indicate that by incorporating the  $M_{dis}$ , a tiny additional number of

Methods	MQuAKE-CF-3K		MQuAKE-T	
	1-edited	All-edited	1-edited	All-edited
MeLLO	28.80	25.30	52.84	53.53
MeLLO w/ $M_{pre}$	29.53	21.20	52.41	54.39
PokeMQA w/o $M_{dis}, M_{gen}$	<b>49.00</b>	<b>29.93</b>	<b>67.99</b>	<b>55.67</b>

Table 3: Results about decoupling effects. Hop-Acc (Metric), GPT-3.5-turbo-instruct (Base model).

Methods	1-edited			All-edited		
	2-hop	3-hop	4-hop	2-hop	3-hop	4-hop
PokeMQA w/o $M_{gen}$	74.00	54.53	46.63	71.20	50.40	36.83
PokeMQA	<b>83.77</b>	<b>64.37</b>	<b>55.13</b>	<b>80.43</b>	<b>61.00</b>	<b>44.70</b>

Table 4: The accuracy of first subquestion from multi-hop questions in MQuAKE-CF-3K on GPT-3.5-turbo-instruct. It should be noted that as long as the first subquestion is answered correctly, the instance is counted as a positive example for the accuracy metric, regardless of whether the language model answered the corresponding multi-hop question correctly or not.

predictions greatly boosts MQA performance under massive edited facts. We conclude that  $M_{dis}$  can greatly enhance the robustness of PokeMQA to large-scale editing with almost no additional computational cost, which is indispensable to maintain the applicability of PokeMQA in real scenarios.



### 5.3 Decoupling Effect Analysis

To better understand the importance of the decoupling operation, we conduct a more fine-grained ablation study to precisely measure its effect. Specifically, we replace the original retriever (contriever-msmarco, Izacard et al., 2021) in MeLLO, which couples the question decomposition with knowledge editing, with pre-detector  $M_{pre}$  and get the new baseline (MeLLO w/  $M_{pre}$ ). The pre-detector  $M_{pre}$  selects the fact statement with highest in-scope probability as the retrieved edited fact. This modification ensures that, in the experiments comparing (PokeMQA w/o  $M_{dis}$ ,  $M_{gen}$ ) and (MeLLO w/  $M_{pre}$ ), coupling remains the sole variable under consideration. The results are shown in Table 3.

It is observed from the experimental results that: I) (PokeMQA w/o  $M_{dis}$ ,  $M_{gen}$ ) has better results across all settings; II) When the edit batch size is 1, it outperforms (MeLLO w/  $M_{pre}$ ) by a substantial margin; III) The retrieval accuracy of the separate pre-detector  $M_{pre}$  does not show significant differences compared to the retriever, so the performance gain comes from decoupling alone, not from other factors like better retrieval accuracy. Thus, we conclude that **decoupling is the key factor contributing to the improved MQA performance**.

### 5.4 Knowledge Prompt Impact Analysis

To elaborate the impact of the knowledge prompt on the editing process more clearly, we calculate the accuracy for the *first subquestion* of multi-hop questions from MQuAKE-CF-3K both with and without the knowledge prompt (PokeMQA w/o  $M_{gen}$ ). The results are shown in Table 4.

Based on the results, it is clear that incorporating knowledge prompts consistently brings an increase in accuracy to the first subquestion across questions with varying hop counts. Thus, we can strongly hypothesize that: I) The impact of knowledge prompts exhibits strong locality, with a significant positive effect on the decomposition and answering of the first subquestion; II) By adding knowledge prompts in context, LLMs can avoid implicit named entity recognition and easily access entity-related facts; III) Although identifying the key entity may not appear challenging for LLMs, it can still pose a distraction to LLMs’ inference ability when embedded in a more complex task, aligning with our insight in Section 3.1.

## 6 Related Work

**Knowledge editing methods.** Knowledge editing focuses on updating factual knowledge to language models and a lot of related research has been carried out. Most of these methods predict updates to the weights of the base model by knowledge locating or meta-learning, and then locally modify parameters (Mitchell et al., 2021, Meng et al., 2022b). Other methods preserve parameters and externally store edit instances (Mitchell et al., 2022, Zhong et al., 2023, Shi et al., 2024). Recent work has identified the limitations of existing editing methods through theoretical analysis (Hase et al., 2023) and performance evaluation (Onoe et al., 2023, Cohen et al., 2024). Our work focuses on addressing one of these challenging tasks: MQA under knowledge editing scenarios.

## 7 Conclusion

In this work, we propose a novel programmable knowledge editing method (PokeMQA) to improve MQA performance and address unreliable reasoning. PokeMQA leverages a scope detector to align LLMs’ behavior with edited facts and incorporates auxiliary knowledge prompt to enrich contextual information. Extensive experiments across three LLMs indicate that PokeMQA helps LLMs answer multi-hop questions in a precise and reliable way.

### Limitations

The limitations of our work are as follows:

- In this work, we demonstrate that our proposed two-stage edited fact retrieval framework exhibits precise retrieval capability and a significant degree of applicability in real-world scenarios. Based on the training paradigm, future work may study how to design a task-specific architecture for the scope detector to advance the retrieval performance.
- Our proposed PokeMQA still struggle to solve complex multi-hop questions with more fact hops. Trying to mitigate the context-following pressure imposed by the overlong context length during the question decomposition process might be a promising research direction.
- Besides, although memory-based editing shows great potential for stable and controlled editing and large-scale editing, the way it

stores edit instances makes it extremely vulnerable to attacks such as memory injection. Therefore, memory-based editing needs to be supported by reliable security technology to reduce its risk in real scenarios.

## Acknowledgements

We thank all reviewers for providing valuable feedback. This work was supported by the National Science and Technology Major Project under Grant No.2023YFF0905400, and the National Natural Science Foundation of China under grants (No.62372211, 62272191), and the International Science and Technology Cooperation Program of Jilin Province (No.20230402076GH, No.20240402067GH), and the Science and Technology Development Program of Jilin Province (No.20220201153GX).

## References

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Huiyuan Chen, Xiaoting Li, Kaixiong Zhou, Xia Hu, Chin-Chia Michael Yeh, Yan Zheng, and Hao Yang. 2022. [Tinykg: Memory-efficient training framework for knowledge graph neural recommender systems](#). In *Proceedings of the 16th ACM Conference on Recommender Systems*, RecSys '22. ACM.
- Jifan Chen, Shih-ting Lin, and Greg Durrett. 2019. Multi-hop question answering via reasoning chains. *arXiv preprint arXiv:1910.02610*.
- Tiejun Chen, Longchao Da, Huixue Zhou, Pingzhi Li, Kaixiong Zhou, Tianlong Chen, and Hua Wei. 2024. Privacy-preserving fine-tuning of large language models through flatness. *arXiv preprint arXiv:2403.04124*.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%\\* chatgpt quality](#).
- Yu-Neng Chuang, Guanchu Wang, Chia-Yuan Chang, Kwei-Herng Lai, Daochen Zha, Ruixiang Tang, Fan Yang, Alfredo Costilla Reyes, Kaixiong Zhou, Xiaoqian Jiang, et al. 2023. Discoverpath: A knowledge refinement and retrieval system for interdisciplinarity on biomedical research. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 5021–5025.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. 2024. Evaluating the ripple effects of knowledge editing in language models. *Transactions of the Association for Computational Linguistics*, 12:283–298.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. *arXiv preprint arXiv:2104.08164*.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghan-deharioun. 2023. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. *arXiv preprint arXiv:2301.04213*.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*.
- Mingyu Jin, Xue Haochen, Zhenting Wang, Boming Kang, Ruosong Ye, Kaixiong Zhou, Mengnan Du, and Yongfeng Zhang. 2024. Prollm: Protein chain-of-thoughts enhanced llm for protein-protein interaction prediction. *bioRxiv*, pages 2024–04.
- Muhammad Khalifa, Lajanugen Logeswaran, Moon-tae Lee, Honglak Lee, and Lu Wang. 2023. Few-shot reranking for multi-hop qa via language model prompting. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15882–15897.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Yunshi Lan, Gaole He, Jinhao Jiang, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. 2021. A survey on complex knowledge base question answering: Methods, challenges and solutions. *arXiv preprint arXiv:2105.11644*.
- Belinda Z Li, Sewon Min, Srinivasan Iyer, Yashar Mehdad, and Wen-tau Yih. 2020. Efficient one-pass end-to-end entity linking for questions. *arXiv preprint arXiv:2010.02413*.
- Gang Liu and Jiabao Guo. 2019. Bidirectional lstm with attention mechanism and convolutional layer for text classification. *Neurocomputing*, 337:325–338.
- Zirui Liu, Guanchu Wang, Shaochen Zhong, Zhaozhuo Xu, Daochen Zha, Ruixiang Tang, Zhimeng Jiang, Kaixiong Zhou, Vipin Chaudhary, Shuai Xu, et al. 2023. Winner-take-all column row sampling for memory efficient adaptation of language model. *arXiv preprint arXiv:2305.15265*.

- Zhibin Lu, Pan Du, and Jian-Yun Nie. 2020. Vgcn-bert: augmenting bert with graph embedding for text classification. In *Advances in Information Retrieval: 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14–17, 2020, Proceedings, Part I 42*, pages 369–382. Springer.
- Vaibhav Mavi, Anubhav Jangra, and Adam Jatowt. 2022. A survey on multi-hop question answering and generation. *arXiv preprint arXiv:2204.09140*.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022b. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2021. Fast model editing at scale. *arXiv preprint arXiv:2110.11309*.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-based model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Yasumasa Onoe, Michael JQ Zhang, Shankar Padmanabhan, Greg Durrett, and Eunsol Choi. 2023. Can lms learn new entities from descriptions? challenges in propagating injected knowledge. *arXiv preprint arXiv:2305.01651*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Van L Parsons. 2014. Stratified sampling. *Wiley StatsRef: Statistics Reference Online*, pages 1–11.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*.
- Dattaraj J Rao, Shraddha S Mane, and Mukta A Paliwal. 2022. Biomedical multi-hop question answering using knowledge graph embeddings and language models. *arXiv preprint arXiv:2211.05351*.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu. 2024. Retrieval-enhanced knowledge editing for multi-hop question answering in language models. *arXiv preprint arXiv:2403.19631*.
- Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitriy Pyrkin, Sergei Popov, and Artem Babenko. 2020. Editable neural networks. *arXiv preprint arXiv:2004.00345*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrusti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85.
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Guozhou Zheng, et al. 2023. Easyedit: An easy-to-use knowledge editing framework for large language models. *arXiv preprint arXiv:2308.07269*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Zhaozhuo Xu, Zirui Liu, Beidi Chen, Yuxin Tang, Jue Wang, Kaixiong Zhou, Xia Hu, and Anshumali Shrivastava. 2023. Compress, then prompt: Improving accuracy-efficiency trade-off of llm inference with transferable prompt. *arXiv preprint arXiv:2305.11186*.
- Yuan Yao, Lorenzo Rosasco, and Andrea Caponnetto. 2007. On early stopping in gradient descent learning. *Constructive Approximation*, 26:289–315.
- Daochen Zha, Kwei-Herng Lai, Kaixiong Zhou, and Xia Hu. 2022. [Towards similarity-aware time-series classification](#). *CoRR*, abs/2201.01413.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023a. Can we edit factual knowledge by in-context learning? *arXiv preprint arXiv:2305.12740*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023b. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. Mquake: Assessing knowledge editing in language models via multi-hop questions. *arXiv preprint arXiv:2305.14795*.

Kaixiong Zhou, Zirui Liu, Rui Chen, Li Li, Soo-Hyun Choi, and Xia Hu. 2022. Table2graph: Transforming tabular data to unified weighted graph. In *IJCAI*, pages 2420–2426.

Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar. 2020. Modifying memories in transformer models. *arXiv preprint arXiv:2012.00363*.

## A Prompt & Supplement Results for PokeMQA

A demonstration example used in prompt is shown in Table 7. The supplement ablation results are shown in Table 8 and Figure 4.

## B Details of Training Dataset Construction

To train our scope detector, including pre-detector  $g_\phi$  and conflict disambiguator  $g_\psi$ , we construct a training dataset  $\mathcal{D}$ . Specifically, we first extract edit triples from MQuAKE-CF and filter out the part sharing the same  $(s, r)$  with fact triples appeared in MQuAKE-CF-3K and MQuAKE-T, constructing a edit dataset  $\mathcal{D}_e = \{e_1, \dots, e_n\}$ . Then we use manually-defined template to convert each edit triple  $e$  into a natural language statement  $s$  and get a edit statement dataset  $\mathcal{D}_{state} = \{t_1, \dots, t_n\}$ . Finally we design a prompt consisting of instruction and demonstrations (shown in Table 9) and prompt Vicuna-13B (Chiang et al., 2023) to generate three diversely phrased atomic questions for each  $t \in \mathcal{D}_{state}$ , building a training dataset  $\mathcal{D} = \{(t_1, q_1), \dots, (t_n, q_n)\}$ . It should be noted that when computing  $SR$  and  $BR$ , or sampling negative instances, instances with the same statement  $t$  should not be considered.

## C Multi-hop Question Answering Dataset Statistics

Table 5 contains the statistics for the two benchmark datasets used in our experiments.

	#Edits	2-hop	3-hop	4-hop	Total
MQuAKE-CF-3K	1	513	356	224	1093
	2	487	334	246	1067
	3	-	310	262	572
	4	-	-	268	268
	All	1000	1000	1000	3000
MQuAKE-T	1 (All)	1421	445	2	1868

Table 5: Statistics of datasets used in experiments

Besides, there are 2786 different edits for MQuAKE-CF-3k and 96 different edits for MQuAKE-T when the size of edit batch is **all**.

## D Details about Scope Detector finetuning

We finetune the pre-detector  $g_\phi$  and conflict disambiguator  $g_\psi$  based on **DistilBERT** (Sanh et al., 2019) and the checkpoint is *distilbert-base-cased* from Huggingface Transformers library (Wolf et al., 2020). We take  $SR + BR - 1$  as the indicator of early stopping.

To fine-tune pre-detector  $g_\phi$ , the learning rate is set as  $1e^{-5}$  with Adam optimizer (Kingma and Ba, 2014), the batch size is set as 1024, and the number of negative samples is 20; To fine-tune conflict disambiguator  $g_\psi$ , the learning rate is set as  $1e^{-5}$  with Adam optimizer, the batch size is set as 256 and the number of negative samples is 1. And the dataset split is 80%/20% for training and validation, without the need of testing.

## E Implementation Details of Baselines

In our experiments, the parameter updating knowledge editing methods, including **FT**, **ROME** and **MEMIT** is implemented by EasyEdit library (Wang et al., 2023). We basically follow the default hyperparameter settings on **LLaMa-2-7B** in library, and make a slight adjustment to ensure the effectiveness of these methods in different experimental settings. We modify the learning rate for ROME and target editing layer for MEMIT, the detailed modifications is shown in Table 6.

Method	MQuAKE-CF-3K			MQuAKE-T	
	1 edited	100 edited	All edited	1 edited	All edited
ROME	$5e^{-1}$	$5e^{-5}$	$5e^{-6}$	$5e^{-1}$	$1.5e^{-1}$
MEMIT	4,5,6,7,8	5,6,7	7	4,5,6,7,8	4,5,6,7,8

Table 6: Detailed hyperparameter modification for ROME and MEMIT.

## F Two-stage edited fact Retrieval

---

**Algorithm 1** Two-stage Edited Fact Retrieval.

---

**Input:** edited memory  $\mathcal{M}$ , pre-detector  $g_\phi$ , conflict disambiguator  $g_\psi$ , atomic question  $q$

- 1: Initialize candidate set  $\mathcal{Z} = \emptyset$
- 2: Initialize final set  $\mathcal{F} = \emptyset$
- 3: */\* Pre-detection stage \*/*
- 4: **for all**  $t_i \in \mathcal{M}$  **do**
- 5:   **if**  $g_\phi(t_i, q) \geq 0.5$  **then**  $\mathcal{Z} = \mathcal{Z} \cup \{t_i\}$
- 6: **end for**
- 7: **if**  $|\mathcal{C}| = 1$  **then**
- 8:   **return**  $t_{i^*}$ , where  $t_{i^*} \in \mathcal{Z}$
- 9: **end if**
- 10: */\* Conflict-disambiguation stage \*/*
- 11: **for all**  $t_i \in \mathcal{C}$  **do**
- 12:   **if**  $g_\psi(t_i, q) \geq 0.5$  **then**  $\mathcal{F} = \mathcal{F} \cup \{t_i\}$
- 13: **end for**
- 14: **if**  $\mathcal{F} \neq \emptyset$  **then**
- 15:   **return**  $t_{i^*}$ , where  $i^* = \arg \max_i g_\psi(t_i, q)$ ,  $t_i \in \mathcal{F}$
- 16: **end if**

---

which can be regarded as an underfitting to question decomposition.

## G Licensing

Vicuna-7B (v1.1) and distilbert-base-cased are released under the Apache License 2.0. LLaMa-2-7B is licensed under the LLAMA 2 Community License. ELQ, ROME, MEMIT, FT are released under the MIT license.

## H Details about Experiments

Because MQuAKE regard a *chain of facts* as an instance and there are three diversely phrased multi-hop questions  $Q$  for each instance, we follow (Zhong et al., 2023), if any of the three questions is considered solved in terms of the specific metric, the instance is considered correct.

The experiments, data, language models in the paper are all in English. We run all experiments on a machine with four NVIDIA A40 GPU. One run of our experiments takes about 15 GPU hours. For all experiments, we use greedy decoding strategy to get the output in text space of language models for reproducibility and report a single run result due to the limited computational resources.

## I Shortcut Reasoning

Given an example from Table 10, although MeLLO appears to successfully combine two facts to arrive at the final answer, its reasoning process does not adhere to the task logic demonstrated in the prompt,

Question: What is the capital city of the country of citizenship of Ivanka Trump’s spouse?  
Entity of Question: **Ivanka Trump, a human.**  
Subquestion: Who is Ivanka Trump’s spouse?  
Generated answer: Ivanka Trump’s spouse is Jared Kushner.  
According to Generated answer, the entity of Subquestion is: Jared Kushner  
Subquestion: What is the country of citizenship of Jared Kushner?  
Generated answer: **Jared Kushner is a citizen of Canada.**  
According to Generated answer, the entity of Subquestion is: Canada  
Subquestion: What is the capital city of Canada?  
Generated answer: The capital city of Canada is Ottawa.  
According to Generated answer, the entity of Subquestion is: Ottawa  
Final answer: Ottawa

Table 7: A in-context demonstration example used in our PokeMQA prompt, here we omit the remaining three demonstrations. **This color** indicate that this part is constructed after being retrieved by knowledge prompt generator from external knowledge base. **This color** indicate that this part is retrieved by scope detector from external memory.

$M_{dis}$ $M_{gen}$	GPT-3.5-turbo-instruct				LLaMa-2-7B				Vicuna-7B				
	MQuAKE-CF-3K		MQuAKE-T		MQuAKE-CF-3K		MQuAKE-T		MQuAKE-CF-3K		MQuAKE-T		
	1 edited	All edited	1 edited	All edited	1 edited	All edited	1 edited	All edited	1 edited	All edited	1 edited	All edited	
-	-	62.60	38.63	76.87	65.63	44.07	28.13	74.68	65.90	44.03	26.27	73.72	66.27
✓	-	62.47	43.23	<b>77.09</b>	<b>78.53</b>	44.00	32.30	74.68	73.82	44.03	29.87	73.72	72.38
-	✓	67.13	40.93	76.93	66.06	<b>44.17</b>	28.17	<b>75.43</b>	66.60	<b>45.90</b>	28.20	<b>74.57</b>	67.29
✓	✓	<b>67.27</b>	<b>45.87</b>	76.98	78.16	44.13	<b>32.83</b>	<b>75.43</b>	<b>74.36</b>	45.83	<b>31.63</b>	<b>74.57</b>	<b>73.07</b>

Table 8: Ablation study results of PokeMQA and its variants in terms of Acc.

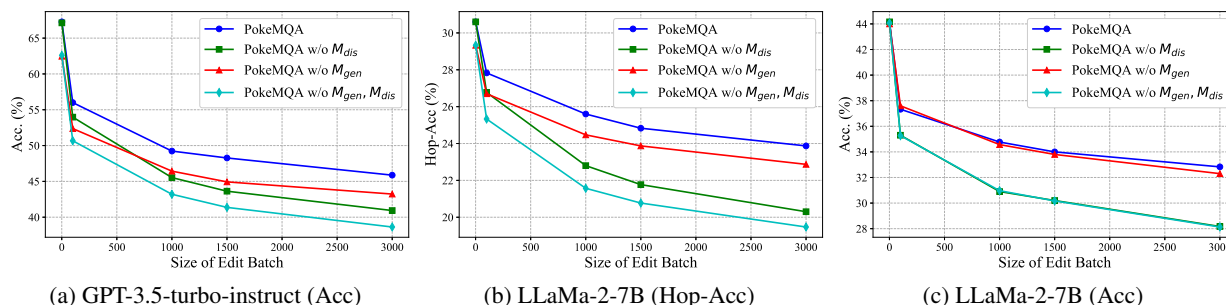


Figure 4: Hop-Acc and Acc across multiple variants of PokeMQA in MQuAKE-CF-3K on GPT-3.5-turbo-instruct and LLaMa-2-7B with edit batches of different sizes.

Please generate three different phrased Questions for each fact

Fact: The univeristy where Bob Dylan was educated is University of Minnesota.  
 Question 1: What is the name of the educational institution where Bob Dylan studied?  
 Question 2: What is the name of the university where Bob Dylan was educated?  
 Question 3: At which university did Bob Dylan receive his education?

Fact: The capital of United Kingdom is Angri.  
 Question 1: What is the name of the capital city of United Kingdom?  
 Question 2: Which city serves as the capital of United Kingdom?  
 Question 3: In which city is the capital of United Kingdom located?

Fact: basketball was created in the country of Spain.  
 Question 1: Where is basketball originated from?  
 Question 2: What is the name of the country of origin of basketball?  
 Question 3: Where did the sport of basketball originate?

Fact: John Coltrane is married to Anne Hathaway.  
 Question 1: Who is John Coltrane’s spouse?  
 Question 2: To whom is John Coltrane married?  
 Question 3: Which person is married to John Coltrane?

Fact: Miranda Kerr is a citizen of Australia.  
 Question 1: Which country is Miranda Kerr a citizen of?  
 Question 2: What is the name of the country that Miranda Kerr belongs to?  
 Question 3: What is Miranda Kerr’s nationality?

Table 9: The rephrased question generation prompt we used for dataset construction. The prompt contains one instruction and some in-context demonstrations.

Question: What is the name of the current head of state of Buckcherry’s country of origin?

Subquestion: What is the country of origin of the band Buckcherry?

Generated answer: Buckcherry’s country of origin is United States of America.

Retrieved fact: The name of the current head of state in United States of America is Joe Biden.

Retrieved fact does not contradict to generated answer, so the intermediate answer is: Joe Biden

Final answer: Joe Biden

Table 10: An example of shortcut reasoning that occurs in the inference process of GPT-3.5-turbo-instruct.