

LLM-based Multi-hop Question Answering with Knowledge Graph Integration in Evolving Environments

Anonymous ACL submission

Abstract

The rapid obsolescence of information in Large Language Models (LLMs) has spurred the development of various techniques for incorporating new facts. To address the ripple effects of altering information, we introduce GMeLLO (Graph Memory-based Editing for Large Language Models), a straightforward yet highly effective method that harnesses the strengths of both LLMs and Knowledge Graphs (KGs). Instead of merely storing edited facts in isolated sentences within an external repository, we utilize established KGs as our foundation and dynamically update them as required. When faced with a query, we employ LLMs to derive an answer based on the relevant edited facts. Additionally, we translate each question into a formal query, tapping into the extensive data within the KG to obtain a more nuanced answer directly from it. In cases of conflicting answers, we prioritize the response derived from the KG as our final result. Our experiments demonstrate a substantial enhancement of GMeLLO over state-of-the-art (SOTA) methods on the MQuAKE benchmark—a dataset specifically designed for multi-hop question answering.

1 Introduction

As the widespread deployment of LLMs continues, the imperative to maintain their knowledge accuracy and up to date, without incurring extensive retraining costs, becomes increasingly evident (Sinitsin et al., 2020). Several approaches have been proposed in prior works to address this challenge, with some focusing on the incremental injection of new facts into language models (Rawat et al., 2020; De Cao et al., 2021; Meng et al., 2022; Mitchell et al., 2022a). Interestingly, certain methodologies in the literature diverge from the conventional path of updating model weights, opting instead for an innovative strategy involving the use of external memory to store the edits (Mitchell et al., 2022b; Zhong et al., 2023). As

Information evolves over time

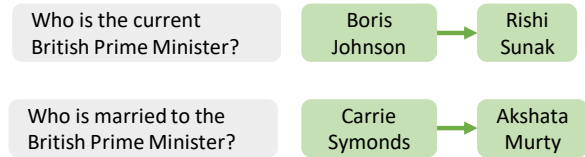


Figure 1: Dynamic nature of information: Changes over time may trigger subsequent modifications. For instance, a transition in the British Prime Minister, such as from Boris Johnson to Rishi Sunak, necessitates corresponding adjustments, like the change in the British Prime Minister’s spouse.

LLMs operate as black boxes, modifying one fact might inadvertently alter another, making it challenging to guarantee accurate revisions. In light of this challenge, opting for an external memory system, rather than directly editing the LLMs, emerges as a prudent choice.

This paper introduces GMeLLO, an effective approach designed to synergize the strengths of LLMs and KGs in addressing the multi-hop question answering task after knowledge editing (Zhong et al., 2023). An illustrative example is presented in Figure 1. Following an update regarding the information of the British Prime Minister, it becomes evident that the corresponding spouse information should also be modified.

As depicted in Figure 2, our GMeLLO comprises the following key steps:

- We utilize LLMs to translate edited fact sentences into triples, employing these triples to update the KG and ensure its information remains up to date.
- Leveraging LLMs again, we analyze a query to extract its relation chain, encompassing the primary entity and its connections with other unknown entities. After populating a template, we convert the relation chain into a formal query and use it to search the updated KG.

- 069 • Based on the query statement, we retrieve
070 the most pertinent edited facts and prompt
071 LLMs to generate an answer in accordance
072 with these facts.
- 073 • In instances where the answer provided by
074 the LLM conflicts with that from the KG, we
075 prioritize the answer from the KG as the final
076 response.

077 LLMs, trained on extensive sentence corpora
078 (Brown et al., 2020; Rae et al., 2022; Chowdhery
079 et al., 2023), are expected to encapsulate a wide
080 range of commonly used sentence structures. As a
081 result, they are invaluable tools for analyzing sen-
082 tences and extracting entities and relations. Once
083 the correct chain of relations and edited triples are
084 obtained, using a formal query to interrogate the
085 KG in a Knowledge-based Question Answering
086 (KBQA) (Cui et al., 2017) manner ensures pre-
087 cision in the retrieval process. In cases where
088 KBQA fails, we still have LLMs for question an-
089 swering (QA) to ensure comprehensive coverage.
090 GMeLLO outperforms current SOTA models on the
091 MQuAKE benchmark, affirming its effectiveness
092 in multi-hop question answering within an evolving
093 environment.

094 2 Related Work

095 The primary focus of this paper lies in exploring
096 enhancing the multi-hop question answering within
097 dynamic scenarios. Therefore, we delve into the
098 related topic of knowledge editing. As highlighted
099 in Yao et al. (2023), two paradigms exist for edit-
100 ing knowledge: modifying model parameters and
101 preserving model parameters.

102 2.1 Modifying Model Parameters

103 In the case of modifying model parameters, this can
104 be further categorized into meta-learning or locate-
105 and-edit approaches. Meta-learning methods, as
106 discussed in (De Cao et al., 2021; Mitchell et al.,
107 2022a), utilize a hyper network to learn the nec-
108 essary adjustments for editing LLMs. The locate-
109 then-edit paradigm, as demonstrated in (Dai et al.,
110 2022; Meng et al., 2022, 2023; Li et al., 2023a;
111 Gupta et al., 2023), involves initially identifying
112 parameters corresponding to specific knowledge
113 and subsequently modifying them through direct
114 updates to the target parameters.

115 2.2 Preserving Model Parameters

116 In the case of preserving model parameters, the
117 introduction of additional parameters or external
118 memory becomes necessary. The paradigm of ad-
119 ditional parameters, as presented in (Dong et al.,
120 2022; Hartvigsen et al., 2022; Huang et al., 2022),
121 incorporates extra trainable parameters into the lan-
122 guage model. These parameters are trained on
123 a modified knowledge dataset, while the original
124 model parameters remain static. On the other hand,
125 memory-based models (Mitchell et al., 2022b;
126 Zhong et al., 2023) explicitly store all edited exam-
127 ples in memory and employ a retriever to extract
128 the relevant edit facts for each new input, guiding
129 the model in generating the edited output.

130 While previous evaluation paradigms have pri-
131 marily focused on validating the recall of edited
132 facts, Zhong et al. (2023) proposed MQuAKE,
133 a benchmark dataset comprising multi-hop ques-
134 tions with either counterfactual edits or temporal
135 edits. This dataset assesses whether methods cor-
136 rectly answer questions where the response should
137 change as a consequence of edited facts. While
138 both GMeLLO and MeLLO (Zhong et al., 2023) are
139 memory-based models targeting multi-hop ques-
140 tion answering in an evolving environment, they
141 differ in the following aspects:

- 142 • MeLLO uses in-context learning to guide
143 LLMs through splitting the question into
144 sub-questions, answering each, and check-
145 ing for contradictions with relevant edit facts.
146 GMeLLO, on the other hand, retrieves a few
147 relevant edit facts for the multi-hop question
148 and presents them along with the question to
149 LLMs for answering
- 150 • Rather than simply storing edited facts as iso-
151 lated sentences in an external memory, we
152 utilize LLMs to translate these sentences into
153 triples and update the KG. Additionally, an-
154 swers are obtained using KBQA to enhance
155 the precision of multi-hop QA within an evol-
156 ving environment.

157 Given that KG is a multi-relational graph consist-
158 ing of entities as nodes and relations among them
159 as typed edges (Saxena et al., 2020), it provides
160 a more straightforward method for representing
161 multi-hop information. Moreover, GMeLLO offers
162 a means to seamlessly integrate the high precision
163 of KBQA (Cui et al., 2017) with the extensive cov-
164 erage of LLMs-based QA, enabling effective multi-

hop question answering in dynamic environments.

3 GMeLLO: Graph Memory-based Editing for Large Language Models

In this section, we explore the details of our method, GMeLLO. Figure 2 provides a visual representation of the GMeLLO framework.

3.1 Utilizing KGs for Storing the Updated Correlated Facts to Enhance Multi-hop Reasoning

KGs play a pivotal role in enhancing the capabilities of LLMs by offering external knowledge for improved inference and interpretability, as demonstrated by recent studies (Pan et al., 2023; Rawte et al., 2023). Apart from merely storing updated information in an external memory, such as a list of separate sentence statements as seen in conventional approaches (Zhong et al., 2023), we utilize the KG to maintain inherent connections and ensure the integration of the latest information.

In our approach, we utilize an off-the-shelf KG, such as Wikidata (Vrandečić and Krötzsch, 2014), as the foundational source. Upon receiving updated facts, we employ LLMs to extract entities and their relationships, forming edited fact triples (Figure 2) that are then used to update the KG.

We incorporate in-context learning (Dong et al., 2023) to ensure the LLMs have a thorough understanding of the task. Furthermore, given the possibility that LLMs may generate relations not present in the KG’s predefined list (Chen et al., 2024), we employ a retriever model to identify the most similar relation from the KG’s list, which is detailed in Section 4.1.6. This relation retrieve procedure is also crucial during relation chain extraction.

3.2 Extracting the Relation Chain of a Question Sentence Using LLMs

With the world changing at a rapid pace, the training data for LLMs can quickly become outdated. Nevertheless, the evolution of patterns tends to occur at a relatively slower pace when compared to the intricate details. In this paper, we employ LLMs to extract the relation chain from a sentence, encompassing the mentioned entity and relations with other unidentified entities. To mitigate varied representations of the same relation, we task LLMs with selecting a relation from a predefined list. Take a question sentence from the MQuAKE dataset as an example,

- Question: What is the capital of the country of citizenship of the child of the creator of Eeyore? 214
- Relation Chain: Eeyore->creator->?x->child->?y->country of citizenship->?z->capital->?id 217

The presented question necessitates a 4-hop reasoning process. With "Eeyore" as the known entity in focus, the journey to the final answer involves identifying its creator, moving on to the creator’s child, obtaining the child’s country of citizenship, and culminating with the retrieval of the country’s capital. The relation chain encapsulates all essential information for arriving at the conclusive answer.

To ensure that LLMs comprehend the task of extracting the relation chain and generate output in a structured template, we employ in-context learning (Dong et al., 2023).

3.3 Converting the Relation Chain into a Formal Query for Retrieving Updated Information from KGs

Once the relation chain is obtained, the next step involves integrating the known entity and the relations into a formal query template. For instance, consider a KG represented in RDF¹ format and a corresponding SPARQL² query. The relation chain elucidated in Section 3.2 should be represented as follows, underscoring the seamless integration of the obtained information into a structured query framework.

```
PREFIX ent: <http://www.kg/entity/> 244
PREFIX rel: <http://www.kg/relation/> 245
SELECT DISTINCT ?id ?label WHERE { 246
  ent:E0 rel:R0 ?x. 247
  ?x rel:R1 ?y. 248
  ?y rel:R2 ?z. 249
  ?z rel:R3 ?id. 250
  ?id rdfs:label ?label. 251
} 252
LIMIT 1 253
```

In this context, "ent" and "rel" serve as prefixes for entity and relation, respectively. The identifier "E0" uniquely represents "Eeyore" within the KG, while the identifiers for "creator," "child," "country of citizenship," and "capital" are denoted as "R0",

¹<https://www.w3.org/RDF/>

²<https://www.w3.org/TR/sparql11-query/>

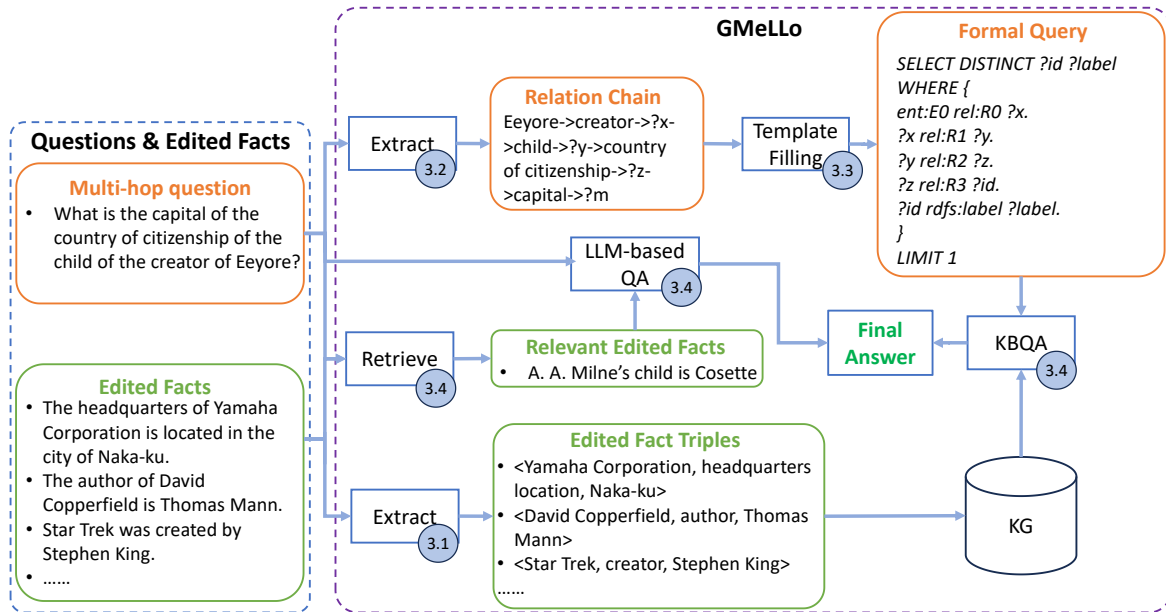


Figure 2: The illustration depicts our proposed method, GMeLLO. We begin by utilizing LLMs to extract entities and relations from edited facts, resulting in a list of edited fact triples. These triples are then used to update a KG. Similarly, we employ LLMs to extract relation chains from a given question. By populating this information into a template, we generate a formal query suitable for use in KBQA (Lan et al., 2022). Simultaneously, we utilize LLMs for question answering, providing an answer based on the relevant edited facts. In cases where the LLM’s answer contradicts that of the KG, we defer to the KG’s answer as the final response.

"R1", "R2", and "R3" respectively. After identifying the entity "?id", we retrieve its string label "?label" as the final answer.

3.4 Enhancing Multi-Hop Question Answering Using Knowledge Graph Integration

When a question arises, we retrieve the "top-x"³ relevant facts using the pretrained Contriever (Izacard et al., 2022) model from a curated list of edited fact sentences. We then prompt the LLMs to generate answers based on the question and these pertinent facts. Compared to the "split-answer-check" pipeline in MeLLO (Zhong et al., 2023), this LLM-based QA method is expected to be simpler and yield more accurate results when the facts are provided accurately. However, addressing multi-hop questions, especially those where the edited facts pertain to intermediary hops, presents a challenge in accurately retrieving the relevant information and performing correct multi-hop question answering. This challenge is particularly pronounced when dealing with a large volume of edited facts. For instance, accurately identifying the relevant fact given the question in Figure 2 and producing

³The "top-x" can be adjusted based on various scenarios. In the majority of cases, it should not exceed 4.

the correct final answer is difficult. To address this issue, we utilize answers from the KG to rectify responses from the LLMs. Once the relation chain and updated triples are derived accurately, the system will yield the correct answer. If the answer is not found within the KG, the system will output nothing, which does not affect the performance of the GMeLLO.

In conclusion, beyond tasking LLMs with question answering, we harness their powerful capabilities for analyzing both edited fact statements and questions. Post-analysis, we convert the edited fact sentences into edited fact triples, subsequently updating the KG. Likewise, we transform the question into a relation chain, culminating in a formal query generated by filling a template, obtaining an answer in a KBQA manner. Our approach leverages KBQA to substitute LLM answers in cases of inconsistency between the two responses. By amalgamating the high precision of KBQA with the expansive coverage of LLMs, our method excels in the multi-hop question answering domain following knowledge editing.

4 Experiment

In the upcoming section, we will conduct experiments to demonstrate the effectiveness of employ-

BaseModel	Method	MQuAKE-CF				MQuAKE-T			
		k=1	k=100	k=1000	k=3000	k=1	k=100	k=500	k=1868
GPT-J-6B	MEMIT	12.3	9.8	8.1	1.8	4.8	1.0	0.2	0.0
GPT-J-6B	MEND	11.5	9.1	4.3	3.5	38.2	17.4	12.7	4.6
GPT-J-6B	MeLLo	20.3	12.5	10.4	9.8	85.9	45.7	33.8	30.7
GPT-J-6B	GMeLLo	50.9	29.2	27.7	27.1	69.9	65.1	64.9	64.8
Vicuna-7B	MeLLo	20.3	11.9	11.0	10.2	84.4	56.3	52.6	51.3
Vicuna-7B	GMeLLo	43.1	20.4	18.1	17.5	75.0	59.0	57.2	57.0

Table 1: Performance results of GMeLLo (ours) on MQuAKE-CF and MQuAKE-T using either GPT-J-6B or Vicuna-7B as the base language model. Following the methodology of [Zhong et al. \(2023\)](#), instances are grouped into batches of size k , where k ranges from 1, 100, 1000, 3000 for MQuAKE-CF, and 1, 100, 500, 1868 for MQuAKE-T. For instance, with the MQuAKE-CF dataset, when $k=100$, the 3000 samples are divided into 30 groups, with the average performance reported as the final result. The metric used is multi-hop accuracy.

ing our GMeLLo methodology.

4.1 Experiment Setup

4.1.1 Dataset

Our experiment centers on the multi-hop question-answering dataset, MQuAKE ([Zhong et al., 2023](#)). This dataset comprises MQuAKE-CF⁴, designed for counterfactual edits, and MQuAKE-T, tailored for temporal knowledge updates. These datasets enable the evaluation of methods under scenarios involving counterfactual changes and real-world temporal updates.

The MQuAKE-CF dataset comprises 3,000 N-hop questions ($N \in \{2, 3, 4\}$), each linked to one or more edits. This dataset functions as a diagnostic tool for examining the effectiveness of knowledge editing methods in handling counterfactual edits. The MQuAKE-T dataset consists of 1,868 instances, each associated with a real-world fact change. Its purpose is to evaluate the efficacy of knowledge editing methods in updating obsolete information with contemporary, factual data.

4.1.2 Baseline Models

To demonstrate the effectiveness of our approach, we conduct comparisons with the following state-of-the-art knowledge editing methodologies.

- MEND ([Mitchell et al., 2022a](#)). It trains a hypernetwork to generate weight updates by transforming raw fine-tuning gradients based on an edited fact.

⁴Our experiments on MQuAKE-CF are carried out on a randomly sampled subset of the complete dataset, comprising 3000 instances (1000 instances for each of 2, 3, 4-hop questions), aligning with the experiments outlined in [Zhong et al. \(2023\)](#).

- MEMIT ([Meng et al., 2023](#)). It updates feed-forward networks across various layers to incorporate all relevant facts.

- MeLLo ([Zhong et al., 2023](#)). It employs a memory-based approach for multi-hop question answering, storing all updated facts in an external memory.

Considering the significant costs associated with training, deploying, and maintaining larger LLMs ([Li et al., 2023b](#)), this paper primarily concentrates on smaller LLMs, specifically GPT-J (6B) ([Wang and Komatsuzaki, 2021](#)) and Vicuna (7B) ([Chiang et al., 2023](#)).

4.1.3 Evaluation Metric

In line with our paper’s central emphasis on multi-hop question answering, we utilize accuracy as the primary metric, to evaluate the methods’ performance in addressing multi-hop inquiries within dynamic environments.

4.1.4 Knowledge Graph Setting

Considering Wikidata’s community-driven nature, guaranteeing a dynamic and comprehensive dataset across a spectrum of knowledge domains, we opt for Wikidata ([Vrandečić and Krötzsch, 2014](#)) as the foundational KG for this experiment. Using LLMs along with 10 <edited fact, edited triple> pairs as samples in the prompt, we extract modified triples from the revised facts with the intention of using them to update the KG. To align the relationships in the questions of test samples with those in Wikidata ([Vrandečić and Krötzsch, 2014](#)), we follow the following steps:

- We select the first 500 item properties⁵ from WikiData as the base relations. Items represent either concrete or abstract entities, such as a person (Piscopo and Simperl, 2019).
- Next, we employ GPT-3.5-Turbo⁶ to examine each multi-hop question in the test samples to determine if it contains any of the base relations.
- Afterward, we rank the frequencies of each relation and choose the top 50 relations as candidates for use in relation chain extraction and edited fact triple extraction.

To stay updated with the latest information on WikiData, we utilize the WikiData API service⁷ and the WikiData Query Service⁸. Since WikiData may contain items with identical labels⁹, we map the entity string in the edited fact triples and the relation chain to WikiData and select the first match as the candidate. We then verify if this entity corresponds to the intended one in the dataset. The correctness of our KBQA result hinges on two crucial criteria:

- The accurate extraction of both edited fact triples and relation chains.
- A precise match between the entity id retrieved from the WikiData API service for each entity string in the edited facts and relation chains and the intended entity id in the dataset.

If the relation chain is found to be incorrect, we conduct an online search on WikiData to determine if the relation chain leads to an entity that could potentially yield an incorrect answer for the specific question, which takes about 1 second.

4.1.5 Prompt Setup and Post-Processing

Compared to MeLLO (Zhong et al., 2023), we adopt a strict evaluation approach, assessing only the first multi-hop question in the MQuAKE datasets for our GMeLLO, instead of considering all three and accepting any one correct. To enhance

⁵<https://www.wikidata.org/w/index.php?title=Special:ListProperties/wikibase-item&limit=500&offset=0>

⁶<https://platform.openai.com/docs/models/gpt-3-5-turbo>

⁷<https://www.wikidata.org/w/api.php>

⁸<https://query.wikidata.org/sparql>

⁹https://www.wikidata.org/wiki/Help:Label/general_principles

the comprehension of the relation chain extraction task by LLMs and ensure outputs adhere to a specified format, we utilize a 3-shot learning approach. This approach entails presenting the model with one 2-hop question sample, one 3-hop question sample, and one 4-hop question sample.

We also implement the in-context learning (Dong et al., 2023) for LLM-based QA. We provide 4 samples in the prompt for MQuAKE-CF: one 1-edit sample, one 2-edit sample, one 3-edit sample, and one 4-edit sample. When $k \geq 100$, we retrieve 4 relevant edit facts for each test sample. When $k=1$, the prompt consists of all the relevant facts for a specific test sample, given that the edit facts in the memory is less than 5.

To address the limitations of GPT-J and Vicuna in conforming to the desired output format, we establish a heuristic rule for extracting essential information from their outputs. For instance, in the context of relation chain extraction, this heuristic is outlined as follows:

- Narrow the attention to the output sentence containing the "->" indicator.
- Divide the sentence based on the "->" delimiter.
- Regard the initial segment as the predicted entity. Subsequently, process the following segments sequentially as relations, provided they do not begin with "?".

4.1.6 Strategies for Managing Unforeseen Relationships

As previously noted, since LLMs may produce relations that are similar in meaning but not identical, we employ the pretrained Contriever model (Izacard et al., 2022) to retrieve the most similar relation (i.e., the closest relation in the embedding space) from the base list of relations. This replacement is performed when undefined relations are encountered during both edited fact triple extraction and relation chain extraction.

4.2 Main Results

As shown in Table 1, our GMeLLO demonstrates significantly superior performance compared to state-of-the-art models on the MQuAKE datasets, including the MQuAKE-CF dataset and MQuAKE-T dataset. Particularly noteworthy is its performance when handling multiple edits simultaneously. When $k=3000$ and using GPT-J as the base

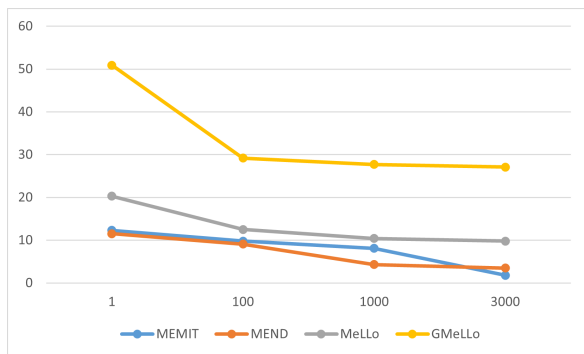


Figure 3: The performance comparison of different methods on MQuAKE-CF dataset when using GPT-J as the base model. The evaluation is conducted with varying numbers of edited instances (k) selected for editing, where k ranges from 1 to 3000.

model, GMeLlo shows an improvement of roughly 18% over MeLlo in MQuAKE-CF, and approximately 30% in MQuAKE-T.

As with many other approaches, we witness a significant decline from $k=1$ to $k=100$. This is understandable, as at $k=1$, all edited facts related to a question are fed into the prompt for LLMs to answer without requiring retrieval. However, the performance stabilizes thereafter. The graph in Figure 3 demonstrates that integrating KBQA enables GMeLlo to maintain higher performance levels, even with an increasing number of edits.

4.3 Ablation Study

To gain a comprehensive understanding of the performance of various components, i.e., LLM-based QA and KBQA, we conduct an experiment to illustrate the impact of LLM-based QA and KBQA as the number of edits increases. As demonstrated in Table 2, the performance of KBQA remains consistent, as all edited facts are converted to triples and all relation chains are extracted from the test questions, regardless of the value of 'k'. However, as the parameter 'k' increases, more edited facts are stored in the external memory. Consequently, selecting the relevant edits and accurately answering the questions becomes increasingly challenging for LLM-based QA.

As depicted in Table 2, when $k=1$ and all relevant facts are provided to the LLMs for question answering, the process proves to be more effective. However, a more realistic scenario involves multiple edits occurring simultaneously, where each question is asked separately (i.e., $k>1$). The performance showcased in this table demonstrates the

effectiveness of our GMeLlo, highlighting that KBQA serves as a valuable enhancement to LLM-based QA within evolving environments.

4.4 The Impact of Entity Ambiguity on WikiData

In Section 4.1.4, we emphasize the importance of not only string matching but also the accurate mapping of entity strings to WikiData, ensuring precision in editing and searching. Table 3 reveals that out of 6015 edited facts in the MQuAKE-CF dataset, 1441 fail to map correctly to the intended entities in WikiData. Within these 1441 inaccurately transformed edit facts, 355 are correct in terms of string matching alone but are erroneously linked to unintended entities. Additionally, out of 3000 questions, the subject entity in 466 questions does not correctly match the intended entities in WikiData. Nearly half of these instances are correctly extracted by LLMs but are mismatched due to entity ambiguity.

We acknowledge that while some entities genuinely share the same labels but represent distinct entities, such as multiple individuals bearing identical names, others are indeed identical entities. This suggests that performance could further improve by addressing this issue and working with an enhanced KG, a direction we leave for future work.

4.5 Error Analysis

Table 2 illustrates that Vicuna exhibits superior performance in directly handling the QA task, particularly when provided with the exact edited facts. Conversely, GPT-J excels in sentence analysis tasks, showcasing its high performance in the KBQA task.

4.5.1 Inferior Performance of GPT-J in QA

Table 2 shows that the performance of GPT-J and Vicuna in conducting QA tasks is comparable on the MQuAKE-CF dataset when $k=1$. However, GPT-J exhibits notably lower performance on the MQuAKE-T dataset under the same conditions. Further analysis revealed that GPT-J struggles in answering questions with only an edited fact pertaining to its intermediary information, such as:

- Facts: The name of the current head of the Philippines government is Bongbong Marcos
- Question: Who is the head of government of the country that Joey de Leon is a citizen of?

BaseModel	Method	MQuAKE-CF				MQuAKE-T			
		k=1	k=100	k=1000	k=3000	k=1	k=100	k=500	k=1868
GPT-J-6B	QA	66.6	11.1	7.2	6.4	20.9	10.7	9.4	9.1
GPT-J-6B	KBQA	24.2	24.2	24.2	24.2	63.5	63.5	63.5	63.5
GPT-J-6B	GMeLLO	50.9	29.2	27.7	27.1	69.9	65.1	64.9	64.8
Vicuna-7B	QA	69.8	15.6	9.1	7.2	94.4	56.9	52.9	52.0
Vicuna-7B	KBQA	14.0	14.0	14.0	14.0	37.3	37.3	37.3	37.3
Vicuna-7B	GMeLLO	43.1	20.4	18.1	17.5	75.0	59.0	57.2	57.0

Table 2: Performance comparison as edited facts increase among various methods. QA involves directly using LLM for answering the multi-hop questions. KBQA involves using LLM to transform edited fact sentences into triples, update WikiData, convert question sentences into relation chains, and generate formal questions for answering in a KBQA manner. GMeLLO combines these methods: opting for QA when KBQA yields no response and choosing KBQA when QA and KBQA answers differ.

Model	Edited Fact	Relation Chain
GPT-J-6B	355/1441	205/466
Vicuna-7B	345/2033	206/317

Table 3: The error rate of entity mapping from entity strings to entities in WikiData. Due to entity ambiguity in WikiData, a single string may correspond to multiple entities. In the context of GPT-J and MQuAKE-CF, '355/1441' in the edited fact indicates that out of 1441 errors in correctly extracting the fact triple, 355 errors stem from entity mapping.

- Predicted Answer: Benigno Aquino III
- Label: Bongbong Marcos

However, it is worth noting that all test samples in MQuAKE-T contain only one edited fact. In contrast, approximately 63.6% of test samples in MQuAKE-CF consist of more than 2 edited facts, which allows GPT-J to connect all the information together, resulting in improved performance.

4.5.2 Inferior Performance of Vicuna in KBQA

After analysis, we discovered that out of the 1868 test samples in the MQuAKE-T dataset, 130 samples did not capture the fact triples correctly due to not adhering to the output format. In addition, only 362 relation chains were accurately returned, whereas GPT-J returned 1382 correct relation chains.

It is important to note that even if the relation chain is incorrect, the KBQA system may still provide the correct answer. For instance, in the case of Vicuna, it consistently returns "citizen->country->head of government". Although this

is mapped to the predefined relation list as "country of citizenship->country->head of government", whereas the golden chain is "country of citizenship->head of government", the predicted path still leads to the correct answer.

In addition, while LLMs consistently identifies relations accurately—such as 'head of state,' 'chief of department,' and 'head of government'—it often makes errors in their sequencing. To address this, we employ Spacy¹⁰ to detect instances where the object of an edited triple is not a person. If it is not, we adjust the sequence of the object and subject in the triple accordingly.

5 Conclusion

In this paper, we present GMeLLO, a method designed for multi-hop question answering in dynamic environments. Except leveraging LLMs for question answering, we also leverage the capabilities of LLMs to extract the triples from edited fact sentence to update KG, and use the capabilities of LLMs to analyze question sentences and generate a relation chain, and finally get the formal query by filling in a formal query template. Finally, we combine KBQA and LLM-based QA to bolster the multi-hop question answering capability within a dynamic environment. This approach capitalizes on the strengths of both LLMs and KGs—leveraging the high coverage of LLMs and the precision of using KGs. By utilizing LLMs for analyzing most question sentences and QA, and KBQA to provide accurate results, we achieve a synergy between the two methodologies.

¹⁰<https://spacy.io/>

593 Limitations

594 Nevertheless, it’s important to note that this inves-
595 tigation is still in its early stages. Although our
596 performance surpasses that of baseline approaches
597 in the multi-hop question answering when editing
598 multiple facts simultaneously, we recognize the po-
599 tential for further improvement. Looking ahead,
600 our future plans involve enhancing GMeLLO in the
601 following key areas:

- 602 • Experiment with more sophisticated prompts,
603 such as Chain of Thought (CoT) (Wei et al.,
604 2022), to elevate performance.
- 605 • Mitigate the entity ambiguity in KGs to fur-
606 ther improve the performance.
- 607 • Pioneering the integration of the strengths in-
608 herent in both LLMs and KGs, we aim to
609 extend their application to diverse research
610 endeavors.

611 References

612 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
613 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
614 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
615 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
616 Gretchen Krueger, Tom Henighan, Rewon Child,
617 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
618 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
619 teusz Litwin, Scott Gray, Benjamin Chess, Jack
620 Clark, Christopher Berner, Sam McCandlish, Alec
621 Radford, Ilya Sutskever, and Dario Amodei. 2020.
622 [Language models are few-shot learners](#). In *Ad-
623 vances in Neural Information Processing Systems*,
624 volume 33, pages 1877–1901. Curran Associates,
625 Inc.

626 Ruirui Chen, Chengwei Qin, Weifeng Jiang, and
627 Dongkyu Choi. 2024. Is a large language model
628 a good annotator for event extraction? In *Proceed-
629 ings of the AAAI Conference on Artificial Intelligence*,
630 volume 38, pages 17772–17780.

631 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
632 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
633 Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion
634 Stoica, and Eric P. Xing. 2023. [Vicuna: An open-
635 source chatbot impressing gpt-4 with 90%* chatgpt
636 quality](#).

637 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
638 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul
639 Barham, Hyung Won Chung, Charles Sutton, Sebas-
640 tian Gehrmann, et al. 2023. Palm: Scaling language
641 modeling with pathways. *Journal of Machine Learn-
642 ing Research*, 24(240):1–113.

Wanyun Cui, Yanghua Xiao, Haixun Wang, Yangqiu
Song, Seung-won Hwang, and Wei Wang. 2017.
[Kbqa: learning question answering over qa cor-
pora and knowledge bases](#). *Proc. VLDB Endow.*,
10(5):565–576.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao
Chang, and Furu Wei. 2022. [Knowledge neurons in
pretrained transformers](#). In *Proceedings of the 60th
Annual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers)*, pages 8493–
8502, Dublin, Ireland. Association for Computational
Linguistics.

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

Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu,
Zhifang Sui, and Lei Li. 2022. [Calibrating factual
knowledge in pretrained language models](#). In *Find-
ings of the Association for Computational Linguistics:
EMNLP 2022*, pages 5937–5947, Abu Dhabi, United
Arab Emirates. Association for Computational Lin-
guistics.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong
Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and
Zhifang Sui. 2023. [A survey on in-context learning](#).

Anshita Gupta, Debanjan Mondal, Akshay Sheshadri,
Wenlong Zhao, Xiang Li, Sarah Wiegrefe, and Niket
Tandon. 2023. [Editing common sense in transfor-
mers](#). In *Proceedings of the 2023 Conference on Empir-
ical Methods in Natural Language Processing*, pages
8214–8232.

Thomas Hartvigsen, Swami Sankaranarayanan, Hamid
Palangi, Yoon Kim, and Marzyeh Ghassemi. 2022.
[Aging with grace: Lifelong model editing with dis-
crete key-value adaptors](#). In *NeurIPS 2022 Workshop
on Robustness in Sequence Modeling*.

Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou,
Wenge Rong, and Zhang Xiong. 2022. [Transformer-
patcher: One mistake worth one neuron](#). In *The
Eleventh International Conference on Learning Rep-
resentations*.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebas-
tian Riedel, Piotr Bojanowski, Armand Joulin, and
Edouard Grave. 2022. [Unsupervised dense informa-
tion retrieval with contrastive learning](#). *Transactions
on Machine Learning Research*.

Yunshi Lan, Gaole He, Jinhao Jiang, Jing Jiang,
Wayne Xin Zhao, and Ji-Rong Wen. 2022. [Complex
knowledge base question answering: A survey](#). *IEEE
Transactions on Knowledge and Data Engineering*.

Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun
Ma, and Jie Yu. 2023a. [Pmet: Precise model editing
in a transformer](#).

699	Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del	language models: Methods, analysis & insights from	758
700	Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023b.	training gopher.	759
701	Textbooks are all you need ii: phi-1.5 technical re-		
702	port.		
703	Kevin Meng, David Bau, Alex Andonian, and Yonatan	Ankit Singh Rawat, Chen Zhu, Daliang Li, Felix Yu,	760
704	Belinkov. 2022. Locating and editing factual asso-	Manzil Zaheer, Sanjiv Kumar, and Srinadh Bhojana-	761
705	ciations in gpt. In <i>Advances in Neural Information</i>	palli. 2020. Modifying memories in transformer	762
706	<i>Processing Systems</i> , volume 35, pages 17359–17372.	models. In <i>International Conference on Machine</i>	763
707	Curran Associates, Inc.	<i>Learning (ICML) 2021.</i>	764
708	Kevin Meng, Arnab Sen Sharma, Alex Andonian,	Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A	765
709	Yonatan Belinkov, and David Bau. 2023. Mass edit-	survey of hallucination in large foundation models.	766
710	ing memory in a transformer. <i>The Eleventh Inter-</i>	<i>arXiv preprint arXiv:2309.05922.</i>	767
711	<i>national Conference on Learning Representations</i>		
712	<i>(ICLR).</i>	Apoorv Saxena, Aditay Tripathi, and Partha Talukdar.	768
713	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea	2020. Improving multi-hop question answering over	769
714	Finn, and Christopher D Manning. 2022a. Fast model	knowledge graphs using knowledge base embeddings.	770
715	editing at scale. In <i>International Conference on</i>	In <i>Proceedings of the 58th Annual Meeting of the As-</i>	771
716	<i>Learning Representations.</i>	<i>sociation for Computational Linguistics</i> , pages 4498–	772
717	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea	4507, Online. Association for Computational Lin-	773
718	Finn, and Christopher D. Manning. 2022b. Memory-	guistics.	774
719	based model editing at scale. In <i>International Con-</i>	Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitry Pyrkin,	775
720	<i>ference on Machine Learning.</i>	Sergei Popov, and Artem Babenko. 2020. Editable	776
721	Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Ji-	neural networks. In <i>International Conference on</i>	777
722	apu Wang, and Xindong Wu. 2023. Unifying large	<i>Learning Representations.</i>	778
723	language models and knowledge graphs: A roadmap.	Denny Vrandečić and Markus Krötzsch. 2014. Wiki-	779
724	<i>arXiv preprint arXiv:2306.08302.</i>	data: a free collaborative knowledgebase. <i>Communi-</i>	780
725	Alessandro Piscopo and Elena Simperl. 2019. What we	<i>cations of the ACM</i> , 57(10):78–85.	781
726	talk about when we talk about wikidata quality: a	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-	782
727	literature survey. In <i>Proceedings of the 15th Interna-</i>	6B: A 6 Billion Parameter Autoregressive Lan-	783
728	<i>national Symposium on Open Collaboration</i> , OpenSym	guage Model. https://github.com/kingoflolz/	784
729	'19, New York, NY, USA. Association for Computing	mesh-transformer-jax.	785
730	Machinery.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	786
731	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie	Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le,	787
732	Millican, Jordan Hoffmann, Francis Song, John	and Denny Zhou. 2022. Chain-of-thought prompt-	788
733	Aslanides, Sarah Henderson, Roman Ring, Susan-	ing elicits reasoning in large language models.	789
734	nah Young, Eliza Rutherford, Tom Hennigan, Ja-	In <i>Advances in Neural Information Processing Systems</i> ,	790
735	cob Menick, Albin Cassirer, Richard Powell, George	volume 35, pages 24824–24837. Curran Associates,	791
736	van den Driessche, Lisa Anne Hendricks, Mari-	Inc.	792
737	beth Rauh, Po-Sen Huang, Amelia Glaese, Joh-	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,	793
738	annes Welbl, Sumanth Dathathri, Saffron Huang,	Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu	794
739	Jonathan Uesato, John Mellor, Irina Higgins, Anto-	Zhang. 2023. Editing large language models: Prob-	795
740	nia Creswell, Nat McAleese, Amy Wu, Erich Elsen,	lems, methods, and opportunities. In <i>Proceedings</i>	796
741	Siddhant Jayakumar, Elena Buchatskaya, David Bud-	<i>of the 2023 Conference on Empirical Methods in</i>	797
742	den, Esme Sutherland, Karen Simonyan, Michela Pa-	<i>Natural Language Processing</i> , pages 10222–10240,	798
743	ganini, Laurent Sifre, Lena Martens, Xiang Lorraine	Singapore. Association for Computational Linguis-	799
744	Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena	tics.	800
745	Gribovskaya, Domenic Donato, Angeliki Lazaridou,	Zexuan Zhong, Zhengxuan Wu, Christopher Manning,	801
746	Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsim-	Christopher Potts, and Danqi Chen. 2023. MQuAKE:	802
747	poukkelli, Nikolai Grigorev, Doug Fritz, Thibault So-	Assessing knowledge editing in language models via	803
748	tiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong,	multi-hop questions. In <i>Proceedings of the 2023</i>	804
749	Daniel Toyama, Cyprien de Masson d’Autume, Yujia	<i>Conference on Empirical Methods in Natural Lan-</i>	805
750	Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin,	<i>guage Processing</i> , pages 15686–15702, Singapore.	806
751	Aidan Clark, Diego de Las Casas, Aurelia Guy,	Association for Computational Linguistics.	807
752	Chris Jones, James Bradbury, Matthew Johnson,		
753	Blake Hechtman, Laura Weidinger, Iason Gabriel,	A Appendix	808
754	William Isaac, Ed Lockhart, Simon Osindero, Laura		
755	Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub,	The appendix will contain further elaboration on	809
756	Jeff Stanway, Lorraine Bennett, Demis Hassabis, Ko-	the details we used.	810
757	ray Kavukcuoglu, and Geoffrey Irving. 2022. Scaling		

A.1 Prompts

The prompts used for edited fact triple extraction, relation chain extraction, and LLM-based QA are depicted in Figures 4, 5, and 6. The edited triple can be regarded as a specialized relation chain, with only one relation between entities and all entities known. All samples in the prompt are selected from the complete MQuAKE-CF dataset, ensuring they are distinct from the test samples.

Prompt for Transforming the Edited Sentences to Triples
Sentence: The headquarters of University of Cambridge is located in the city of Washington, D.C.
Relation Chain: University of Cambridge->headquarters location->Washington, D.C.
.....
Given the above samples, please help me analyze the relation chain of the following sentence. All the relations should be selected from ['country of origin', 'sport', ...].
Sentence: The chief executive officer of Boeing is Marc Benioff
Relation Chain:

Figure 4: The prompt used for transforming edited fact sentences to triples.

Prompt for Transforming the Question Sentences to Relation Chains
Question: What is the birthplace of the author of "The Little Match Girl"?
Relation Chain: The Little Match Girl->author->x->place of birth->y
.....
Given the above samples, please help me analyze the relation chain of the following sentence. All the relations should be selected from ['country of origin', 'sport', ...].
Question: What is the continent where the CEO responsible for developing Windows 8.1 was born?
Relation Chain:

Figure 5: The prompt used for transforming question sentences to relation chains.

Prompt for LLM-based QA
Facts: Hans Christian Andersen was born in the city of Brittany
Question: What is the birthplace of the author of "The Little Match Girl"?
Answer: Brittany
.....
Facts: Windows 8.1 was developed by Boeing; The chief executive officer of Boeing is Marc Benioff; California is located in the continent of Europe; Marc Benioff was born in the city of California
Question: What is the continent where the CEO responsible for developing Windows 8.1 was born?
Answer:

Figure 6: The prompt used in LLM-based QA.

A.2 Relations

After filtering by GPT-3.5-Turbo, the first 50 relations utilized in MQuAKE-CF dataset are: ['country of origin', 'sport', 'country of citizenship', 'capital', 'continent', 'official language',

'head of state', 'head of government', 'creator', 'country', 'author', 'headquarters location', 'place of birth', 'spouse', 'director / manager', 'religion or worldview', 'genre', 'work location', 'performer', 'manufacturer', 'developer', 'place of death', 'employer', 'educated at', 'member of sports team', 'head coach', 'languages spoken, written or signed', 'notable work', 'child', 'founded by', 'location', 'chief executive officer', 'original broadcaster', 'chairperson', 'occupation', 'position played on team / speciality', 'member of', 'language of work or name', 'director', 'league', 'home venue', 'native language', 'composer', 'place of origin (Switzerland)', 'officeholder', 'religious order', 'publisher', 'original language of film or TV show', 'ethnic group', 'military branch'].

After GPT-3.5-Turbo filtering, the MQuAKE-T dataset includes a total of 35 relations. The relation list is ['head of government', 'country of citizenship', 'head of state', 'country of origin', 'country', 'headquarters location', 'location', 'sport', 'performer', 'genre', 'developer', 'employer', 'manufacturer', 'place of death', 'place of birth', 'author', 'member of', 'capital', 'member of sports team', 'chief executive officer', 'notable work', 'director / manager', 'original broadcaster', 'creator', 'work location', 'educated at', 'located in the administrative territorial entity', 'head coach', 'place of publication', 'location of formation', 'director', 'producer', 'transport network', 'continent', 'child']

A.3 Further Details

All experiments are conducted on NVIDIA RTX A5000 GPUs, with the temperature of LLMs set to 0 across all tasks.