

# Retrieval and Reasoning on KGs: Integrate Knowledge Graphs into Large Language Models for Complex Question Answering

Anonymous ACL submission

## Abstract

Despite Large Language Models (LLMs) have performed impressively in various Natural Language Processing (NLP) tasks, their inherent hallucination phenomena severely challenge their credibility in complex reasoning. Combining explainable Knowledge Graphs (KGs) with LLMs is a promising path to address this issue. However, structured KGs are difficult to utilize, and how to make LLMs understand and incorporate them is a challenging topic. We thereby reorganize a more efficient structure of KGs, while designing the KG-related instruction tuning and continual pre-training strategies to enable LLMs to learn and internalize this form of representation effectively. Moreover, we construct subgraphs to further enhance retrieval capabilities of KGs via CoT reasoning. Extensive experiments on two KGQA datasets demonstrate that our model achieve convincing performance compared to strong baselines<sup>1</sup>.

## 1 Introduction

The emergence of Large Language Models (LLMs) (OpenAI, 2022, 2023; Bubeck et al., 2023; Yang et al., 2023) has attracted widespread attention over the recent years. They demonstrates remarkable reasoning capabilities, managing to solve complex problems through step-by-step thinking and planning (Wei et al., 2022; Khot et al., 2023). However, the reasoning of LLMs is not always reliable and may conflict with factual reality, known as hallucination (Wang et al., 2023; Huang et al., 2023; Zhang et al., 2023). This limitation will restrict the application of LLMs in fields that require high reliability, such as healthcare, law and science. Knowledge Graphs (KGs) store high-quality common sense or domain-specific knowledge in structured triplets. Due to their reliability and interpretability, integrating KGs into LLMs is considered a promising approach to alleviate hallucinations of LLM

reasoning (Pan et al., 2024). Therefore, researchers have never ceased their attempts to integrate KGs with language models (Zhang et al., 2019; Liu et al., 2020; Lewis et al., 2020; Sun et al., 2021), with Knowledge Graph Question Answering (KGQA) being one critical task among them.

The KGQA task faces two challenges: 1) one is how to accurately retrieve specific knowledge from KGs, and 2) the other is to enable reasoning model to understand and utilize this structured knowledge. For the first challenge, some research (Sun et al., 2019; Baek et al., 2023; Jiang et al., 2023b) adopt a direct retrieval approach, using the question as a query and the triples from the KGs as retrieval candidates, employing sparse or dense retrieval techniques to identify the most relevant candidates with the query. However, this way makes it difficult to model the semantic relevance between structured triples and unstructured queries. Besides, the triples that are semantically weakly relevant to the queries may instead be important intermediate knowledge, especially in multi-hop question answering. Another research (Sun et al., 2020; Lan and Jiang, 2020; Gu and Su, 2022; Ye et al., 2022; Yu et al., 2023) transform the question into an executable structured query statement (e.g., SPARQL) and performs the query retrieval in KGs. But it exists the problem of the generating query that are non-executable or executed incorrectly (Yu et al., 2023). For the latter challenge, since LLMs are primarily pre-trained on unstructured text, they may not effectively comprehend and utilize knowledge in the structured form. Consequently, existing methods usually convert KG content to natural language (He et al., 2024; Ye et al., 2024) or linearized triplets (Luo et al., 2024). Nevertheless, natural language form adds redundant tokens, while linearized representation disrupts the structural information inherent within the KG.

To address the above problems, this paper first introduces a novel subgraph-based retrieval-

<sup>1</sup>All the code, data and models will be publicly available at <https://anonymous.com>

augmented method. Specifically, we construct a series of subgraphs via Chain-of-Thought (CoT) (Wei et al., 2022), where subgraphs enrich the semantic information of the candidate knowledge, and CoT offers intermediate reasoning steps involved in multi-hop question answering, aiding the retrieval model in recalling useful intermediary KG knowledge. We then design an efficient KG representation using YAML format to reduce input redundancy, and this organization method does not disrupt the intrinsic structure within the KG. Additionally, we propose three KG-level tasks (including entity, relationship and graph) for instruction tuning and pre-training of KG data to enhance LLM’s understanding of KGs. To further strengthen the reasoning capabilities of LLMs utilizing KGs, we generate explicit reasoning process data with larger open-source LLMs and train our reasoning models with these synthetic datasets.

In summary, our contributions are as follows:

- We introduce a novel and efficient representation for KGs, the YAML format, which reduces token redundancy by approximately 25% compared to the traditional triple format. Combined with our proposed KG-related task tuning, LLMs are able to comprehend and leverage KGs in YAML format to accomplish complex reasoning tasks.
- We integrate the reasoning process and subgraph into knowledge retrieval, which aids in recalling useful intermediate knowledge for reasoning.
- In our experiments conducted on LLaMA2-7b-Chat, our approach has been validated on two challenging KGQA datasets, achieving promising performance compared to strong baselines. Further experimental analysis indicates the generalizability to other LLMs as well.

## 2 Related Work

The KGQA task enables models to answer questions by integrating common sense or domain-specific knowledge from KGs. Current approaches to KGQA can be categorized into three types: embedding-based, semantic parsing-based and retrieval-augmented. Embedding-based methods project entities and relations from KGs into an embedding space, and utilize key-value memory networks (Miller et al., 2016), sequence modeling (He et al., 2021), or graph neural networks (Yasunaga et al., 2021) to learn the reasoning process between questions and the entities and relations. Semantic parsing-based methods utilize the semantic parsing

model to convert questions into structured query language oriented towards the knowledge base (e.g. SPARQL), and then execute it to search answers from the KGs (Sun et al., 2020; Lan and Jiang, 2020; Gu and Su, 2022; Ye et al., 2022; Yu et al., 2023). However, semantic parsing-based methods rely on retrieving answers from knowledge bases, overlooking the reasoning capabilities of models. Retrieval-augmented methods combine KGs with the intrinsic reasoning capabilities of models. They first retrieve question-relevant knowledge triples or subgraphs from the KGs, and then leverage this retrieved knowledge to enhance the factualness of the reasoning. Sun et al. (2018) propose the GraftNet which utilizes entity linking to retrieve subgraphs. Subsequently, many works adopt effective dense retrieval models as their retrieval modules, such as PullNet (Sun et al., 2019), SR (Zhang et al., 2022), DiFar (Baek et al., 2023), UniKGQA (Jiang et al., 2023b), etc. Today, NLP has entered the era of LLMs, where Retrieval-Augmented Generation (RAG) enables these models to effectively leverage external knowledge to accomplish various tasks (Lewis et al., 2020; Gao et al., 2024; Kim et al., 2023; Li et al., 2023a). Wang et al. (2023) retrieve knowledge from KGs to verify and correct the factual within CoT, resulting in the generation of more precision responses. Yu et al. (2023) utilize a larger-scale retriever to enhance retrieval performance and generate both semantic parsing expressions and inference results in the generation phase, compensating for their respective shortcomings by integrating the two approaches.

## 3 Methodology

In this section, we present our proposed KGQA method, which leverages a subgraph-based retrieval-augmentation generation paradigm. First, we introduce the overall inference process of our method, including the KG retrieval module and the KG reasoning module. Then, we detail the training processes for the two modules.

### 3.1 Overview

As Figure 1 shows, our KGQA method includes two modules: KG retrieval model and KG reasoning LLM. Given a question  $q$  and a knowledge graph  $\mathcal{G} = \{t_i\}_i^n$ , where  $t_i = (e_h^i, r^i, e_t^i) \in \mathcal{E} \times \mathcal{R} \times \mathcal{E}$  is a knowledge triple;  $\mathcal{E}, \mathcal{R}$  are the set of entities and relationships;  $e_h, r, e_t$  are the head entity, relationship and tail entity, respectively. Af-

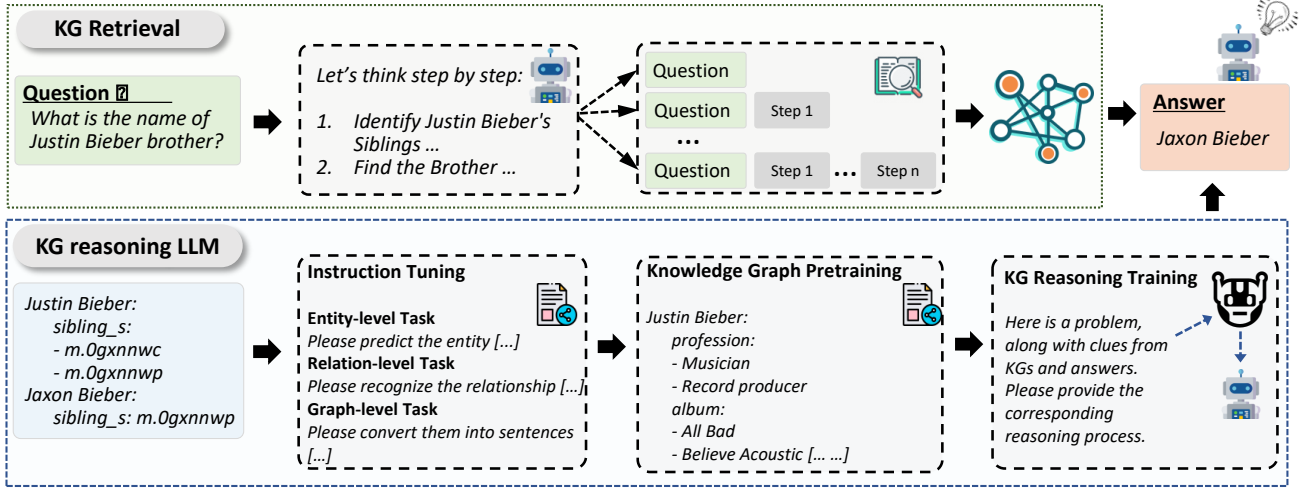


Figure 1: Illustration of our KGQA method. It contains two modules, KG Retrieval Model and KG Reasoning LLM.

ter we train the KG retrieval model  $R_\phi$  and the KG reasoning LLM  $\mathcal{M}_\theta$ , in the inference stage, the LLM  $\mathcal{M}_\theta$  first plans the problem and generates a reasoning process with CoT prompting:

$$\{c^1, \dots, c^j\} = \mathcal{M}_\theta(p_{cot} \oplus q), \quad (1)$$

where  $c^j$  is the  $j$ -th step reasoning process and  $p_{cot}$  is the CoT prompting as shown in Prompt 1,  $\oplus$  means the concatenation operator.

#### Prompt 1: Generating CoT for Retrieval

Please think step by step and then answer the given question.

##### Here are some examples:

**Input:** <Demonstration Question>  
**CoT:** Let's think step by step. <Demonstration CoT>  
**### Output:** <Demonstration Answer>

**Input:** <Question>  
**CoT:** Let's think step by step.

Then, we progressively concatenate the reasoning process with the question as queries to retrieve knowledge:  $q^j = q \oplus c^1 \oplus \dots \oplus c^j$  ( $q^0 = q$ ). For each candidate knowledge triple  $t$ , we integrate the surrounding subgraph information  $\mathcal{G}_t = \{(e_h, r, e_t) | e_h = e_h^t \vee e_t = e_t^t\}$ . The retrieval can be formalized as follows:

$$\mathcal{T} = \text{Top}_k \sum_j f(R_\phi(q^j), R_\phi(t \oplus \mathcal{G}_t)), \quad (2)$$

where  $f$  is the similarity function between the query representation and the candidate representation (e.g. cosine similarity or dot-product similarity),  $\mathcal{T}$  is the set of top- $k$  candidates retrieved that are most relevant to the query.

After retrieval, the candidate set is transformed into YAML format and serves as part of the input for the KG reasoning LLM, which reasons and outputs the final answer through Prompt 2.

#### Prompt 2: Utilizing KG to Reason

Please think step by step and then answer the given question. Please keep the answer as simple as possible and return all the possible answers as a list. If there are hints, please combine this information to answer.

##### Here are some examples:

**Input:** <Demonstration Question>  
**Hints:** <Demonstration Knowledge Graph>  
**CoT:** Let's think step by step. <Demonstration CoT>  
**### Output:** <Demonstration Answer>

**Input:** <Question>  
**Hints:** <Knowledge Graph>  
**CoT:** Let's think step by step.

### 3.2 Subgraph-based Retrieval via CoT

Retrieving relevant and useful knowledge from KGs is critical for the KGQA tasks. Benefiting from the increasingly advanced dense retrieval, we can obtain relevant knowledge through direct retrieval, without the need for elaborate techniques such as semantic parsing and entity linking (Baek et al., 2023). However, the semantic expression of individual knowledge in KGs is limited, and the semantic relationship between knowledge and questions is not directly related in multi-hop question answering. Therefore, we consider incorporating neighboring knowledge information and reasoning processes when retrieving knowledge.

We employ the contrastive learning to train our

retrieval model, the training loss is:

$$\mathcal{L} = -\log \frac{\exp(f(R_\phi(q^j), R_\phi(t^+ \oplus \mathcal{G}_{t^+})))}{\sum_{t \in \tau} \exp(f(R_\phi(q^j), R_\phi(t \oplus \mathcal{G}_t)))}, \quad (3)$$

where  $\tau$  contains all triplets in the same batch,  $t^+$  is the positive sample and others are negative samples. In our method, we take all the knowledge triples on the path from the question entity to the answer entity as positive samples, and randomly sample from the remaining triples as negative samples.

Different from the inference stage, we only use the LLaMA2-7b-Chat model, which has not been specifically trained for KG tasks, to generate the reasoning process for training. This method allows for the complete decoupling of the training of the retrieval and reasoning models, enabling them to be trained independently and in parallel. To address the inconsistency in CoT quality during training and inference, we employ rationalization prompting (Prompt 3<sup>2</sup>) during training, providing the answer in the prompt so that the LLM can generate a reasonable reasoning process based on the answer.

### Prompt 3: Generating CoT for Training

Here is a problem, along with (clues from a knowledge graph and) the answer. Please provide the corresponding reasoning process.

**Here are some examples:**

**Input:** <Demonstration Question>  
**(Clues:** <Demonstration Knowledge Graph>  
**Answer:** <Demonstration Answer>  
**### Output:** <Demonstration CoT>

**Input:** <Question>  
**(Clues:** <Knowledge Triples>  
**Answer:** <Answer>  
**### Output:**

### 3.3 Utilizing KGs Effectively and Efficiently in LLMs

KGs are essentially structured knowledge, while LLMs are typically pre-trained on unstructured text. To bridge this gap and enable LLMs to better understand and utilize the structured knowledge, we propose a simplified representation for KGs. Additionally, we employ instruction tuning and continual pre-training to ensure that LLMs internalize both the knowledge.

<sup>2</sup>Prompt 3 applies to both retrieval training and reasoning training, and KG information is only provided during reasoning training (in section 3.4).



Figure 2: An example of triple and YAML format KG.

**YAML Format KG.** In general, the retrieved knowledge triples may exhibit many literal similarities, such as having the same head entity or relation across multiple triples. If we linearize these triples directly as input for the reasoning LLM, it will result in significant token redundancy, thereby impacting the efficiency of the model’s inference. Therefore, we try to represent the KG in a more efficient format. Our approach uses the YAML format, a data serialization language with a simple syntax. As shown in Figure 2, YAML uses indentation to represent hierarchical relationships. We treat different head entities as the first-level relationship, different relationships under the same head entity as the second level, and different tail entities under the same head entity and relationship as the final level.

**KG-oriented Instruction.** For general-purpose LLMs, representing KGs in YAML format is unfamiliar and infrequently encountered in their pre-training corpora. Therefore, to enable LLMs to understand KGs in YAML, we design three types of graph-related instruction-tuning tasks: 1) **Entity-level tasks**, where the LLM is required to reason the entity according to neighbors; 2) **Relationship-level tasks**, where the task is to reason the relationship between entities; 3) **Graph-level tasks**, where the LLM needs to understand the semantic of KGs and converts to natural language. As shown in Table 5, we design three different instructions for each type of task and denote the instruction prompt as  $\mathcal{I}$ . For entity-level and relationship-level instruction tasks, we automatically construct them based on the data in the KG without the need for additional manual annotation. For graph-level instruction tasks, we utilize existing high-quality KG-to-text datasets (Gardent et al., 2017). The training loss of KG instruction is:

$$\mathcal{L}_{instruct} = -\sum_l^L y^l \log p(\hat{y}^l | \mathcal{I}(x), y^{<l}), \quad (4)$$

where  $(x, y)$  is the input-output pair,  $L$  is the length of  $y$ ,  $y^l$  is the  $l$ -th token,  $y^{<l}$  means tokens before  $l$ -th token,  $\hat{y}^l$  is the predicted  $l$ -th token.

**Continual KG Pre-training.** To further learn the structured knowledge embedded in KGs, we propose the continual KG pre-training method. We serialize the entire KG in YAML format and train it by the next token prediction:

$$\mathcal{L}_{pretrain} = - \sum_l^L x^l \log p(\hat{x}^l | x^{<l}), \quad (5)$$

where  $x$  is the pre-training data.

### 3.4 KG-based Reasoning Training

In section 3.3, we enhance the LLM’s understanding of the specialized structured representation of KG, without explicitly teaching the LLM to use KG for reasoning. In practical scenarios, we need to address two issues: 1) How to utilize KG for multi-hop reasoning; 2) How to manage the retrieved noisy knowledge that lacks crucial task-related information or contains irrelevant redundant information. To address these issues, we use a retrieval model that has not been fine-tuned for KGQA tasks to retrieve noisy knowledge, and a more powerful LLM to generate high-quality reasoning processes for questions based on retrieved knowledge and answers with Prompt 3. After obtain the knowledge and reasoning processes, we train our reasoning LLM with the loss function defined in Equation 4.

## 4 Experiments

### 4.1 Experimental Settings

**Datasets and Evaluation Metrics.** To evaluate the effectiveness of our proposed KGQA method, we conduct experiments on two popular and challenging datasets: WebQSP (Yih et al., 2015) and CWQ (Talmor and Berant, 2018). Both two datasets are created from the Freebase KG (Bollacker et al., 2008). We report more details in Table 4. Following previous work (Jiang et al., 2023b), we take the Hits@1 and F1 as evaluation metrics for WebQSP and CWQ. Hits@1 is a metric for measuring the accuracy of the Top-1 answer. For generative tasks, the order of generation does not imply the probability of the answers. Therefore, we treat all generated responses as the Top-1 answer. Given a question may have multiple answers, F1 balances precision and recall of the

predicted answers, and is used to assess the overall coverage of the model’s predictions.

**Implementation Details.** In our main experiments, we take LLaMA2-7b-Chat<sup>3</sup> as the reasoning backbone model and BGE-1.5-en-base<sup>4</sup> as the retrieval backbone model. We finetune the retrieval model on the training set of WebQSP and CWQ for 5 epochs. The learning rate is set to 1e-5 and the batch size is set to 64. We search for a path in Freebase that starts with a question entity and ends with an answer entity (limiting the length of the path to no more than 5), treating all entities in the path as positive samples of the query, and randomly sampling 6 triples as negative samples. We construct 270k entity-level and 540k relationship-level instruction data from Freebase, and the WebNLG dataset (Gardent et al., 2017) as graph-level instruction data. We tune the reasoning model for 2 epochs with the learning rate set to 2e-6 and batch size set to 64. Then, we perform continual pre-training on the Freebase data using the same setting. For KG-based reasoning training, we use the WebQSP and CWQ training sets as queries to retrieve knowledge from KG using BGE-1.5-en-base. Then, we employ LLaMA2-70b-Chat<sup>5</sup> to generate high-quality reasoning processes, which are used to train our reasoning model. The training is conducted for 5 epochs with the learning rate set to 2e-6 and batch size set to 64. In the inference stage, the first 3 samples from the WebQSP training set are added as demonstrations before each question. For each question, we use our retriever to retrieve the top-20 triples most relevant to it. For generation, we adopt top- $p$  sampling with the temperature set to 0.85 and  $p$  set to 0.9, and the generation length is 512 tokens. To enhance inference speed, model inference is based on the vLLM library (Kwon et al., 2023).

### 4.2 Baselines

We compare our method with the following competitive KBQA baselines:

- NSM (He et al., 2021) proposes a teacher-student framework where the teacher model learns supervision signals for intermediate reasoning processes through forward and backward reasoning, which are then conveyed to the student model for multi-hop inference.

<sup>3</sup><https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>

<sup>4</sup><https://huggingface.co/BAAI/bge-base-en-v1.5>

<sup>5</sup><https://huggingface.co/meta-llama/Llama-2-70b-chat-hf>

- **Transfernet** (Shi et al., 2021) utilizes the graph attention mechanism to capture the relevance among questions, entities, and relationships, guiding a step-by-step traversal on the KG towards the answer.
- **SR + NSM (+E2E)** (Zhang et al., 2022) proposes a effective subgraph retriever to retrieve the most relevant relation-path for reasoning and then utilizes the NSM to reason. **E2E** denotes further jointly finetuning the SR + NSM.
- **QGG** (Lan and Jiang, 2020) is a semantic parsing based approach that incorporates constraints and extends relational paths in the process of generating query graphs.
- **UniKGQA** (Jiang et al., 2023b) unifies the retriever and reasoning module into a single model.
- **DECAF** (Yu et al., 2023) proposes a method for joint generating semantic parsing forms and direct answers, significantly improving the executability of semantic parsing forms.
- **StructGPT** (Jiang et al., 2023a) utilizes LLMs’ tool-using capabilities to interactive between LLMs and KGs, which facilitates multi-hop reasoning through iterative interactions.
- **KD-CoT** (Wang et al., 2023) retrieves relevant knowledge from the KG during the reasoning process, progressively verifying and correcting facts in the reasoning process.
- **RoG** (Luo et al., 2024) leverages the powerful generative and planning capabilities of LLMs to generate reasoning paths. It retrieves corresponding knowledge from knowledge graphs based on these paths and synthesizes various reasoning paths to deduce the final answer.

### 4.3 Main Results

Table 1 shows the results of our model and other baselines on WebQSP and CWQ. Firstly, general-purpose LLMs do not perform well on KGQA tasks, with neither LLaMA2-7b-Chat nor the ChatGPT able to match the performance of KGQA-specific models, especially in the more challenging CWQ dataset. This means that LLMs still have significant room for improvement in their ability to understand and utilize structured knowledge graphs for complex reasoning. Our approach improves Hits@1 by 15-20% compared to these strong general-purpose LLMs. Currently, the state-of-the-art (SOTA) models for KGQA are RoG and DECAF, which are based on retrieval-augmentation and semantic parsing respectively, with backbone models that have over a billion parameters. In terms

Models	WebQSP		CWQ	
	Hits@1	F1	Hits@1	F1
NSM	68.7	62.8	47.6	42.4
TransferNet	71.4	-	48.6	-
SR + NSM	68.9	64.1	50.2	47.1
SR + NSM + E2E	69.5	64.1	49.3	46.3
QGG	73.0	73.8	36.9	37.4
UniKGQA	77.2	72.2	51.2	49.0
DECAF	82.1	<b>78.8</b>	-	-
LLaMA2-7b-Chat	59.5	34.0	34.0	22.7
StructGPT	69.6	-	-	-
ChatGPT	75.6	-	48.9	-
ToG + GPT4	82.6	-	69.5	-
KD-CoT	68.6	52.5	55.7	-
RoG	<u>85.7</u>	70.8	<u>62.6</u>	<b>56.2</b>
Ours	<b>91.5</b>	<u>74.0</u>	<b>68.7</b>	<u>55.6</u>

Table 1: Performance of our model and different baselines on two KGQA datasets. **Bold** and underline represent the best and the second best result, respectively.

Models	Hits@1	Precision	Recall	F1
Ours	68.7	<b>56.4</b>	63.0	<b>55.6</b>
w/o. SubKG-R	65.4	52.9	59.7	52.2
w/o. CoT-R	66.1	52.8	60.3	52.5
w/o. KG-IT	68.0	55.8	62.3	55.1
w/o. KG-PT	<b>69.4</b>	53.8	<b>63.9</b>	54.1
w/o. KG-RT	42.6	34.0	37.0	32.3

Table 2: Ablation study on CWQ. **R**, **IT**, **PT** and **RT** denote retrieval, instruction tuning, continual pretraining and reasoning training, respectively.

of the Hits@1 metric, our method comprehensively surpasses the existing SOTA, especially in the WebQSP dataset, where we achieve a breakthrough of more than 90% for the first time. Compared to RoG, our method shows a significant improvement in 6% Hits@1 on both WebQSP and CWQ. Overall, our method is comparable to the SOTA models in terms of the F1 score. On WebQSP, it falls short of DECAF but outperforms RoG by 3%, and on CWQ, it is on par with RoG.

## 5 Analysis and Discussion

### 5.1 Ablation Study

We conduct ablation experiments on CWQ to analyze the contributions of KG retrieval module and KG reasoning module. As shown in the experimental results in Table 2, each module in our method is indispensable. The most crucial component is KG reasoning training; without it, the model’s perfor-

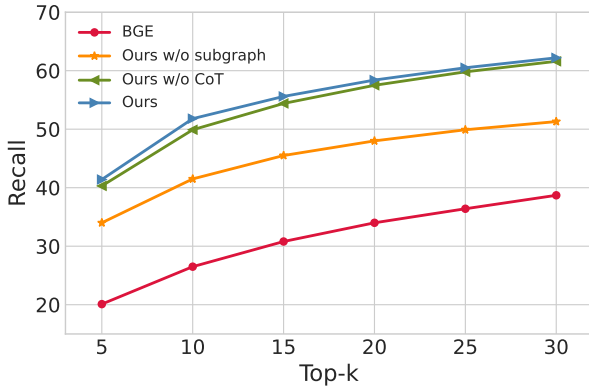


Figure 3: Comparison of recall ability of different retrieval models.

456 mance plummets from 68.7% to 42.6% in Hits@1.  
 457 This indicates that even if LLMs encode KG infor-  
 458 mation and understand its semantics, it is in vain if  
 459 LLMs fail to utilize KG for reasoning. The second  
 460 key component is the retrieval module. Experi-  
 461 ments show that the roles of subgraph information  
 462 and the reasoning process are complementary, and  
 463 their combined use maximizes effectiveness. Lack-  
 464 ing either can lead to a 3% reduction in the model’s  
 465 performance. Compared to the reasoning process,  
 466 subgraph information is more crucial, indicating  
 467 that effectively encoding the semantic information  
 468 of KG in the retrieval model remains the key issue.  
 469 Finally, command fine-tuning and continued  
 470 pre-training also have a positive impact on model  
 471 performance. Instruction tuning can improve the  
 472 model’s performance by about 0.7% across all met-  
 473 rics. Continued pre-training enhances the model’s  
 474 understanding of KG semantics, which helps to  
 475 filter out irrelevant knowledge, thereby improving  
 476 the model’s precision and F1 score.

## 477 5.2 Retrieval Evaluation

478 The performance of retrieval-augmented KGQA  
 479 models is largely dependent on the quality of the  
 480 retrieval process (Jiang et al., 2023b). We expect  
 481 retrieval models to exhibit exceptional recall capa-  
 482 bilities to cover as much useful intermediate knowl-  
 483 edge as possible. This is because while reasoning  
 484 LLMs may learn to filter out irrelevant information  
 485 through training, they struggle to compensate for  
 486 the absence of crucial information. Therefore, we  
 487 compare the recall ability of *our* retrieval model,  
 488 *ours w/o subgraph*, *ours w/o CoT*, and the *BGE*  
 489 model (results are shown in Figure 3). It is evi-  
 490 dent that our retrieval model has a higher recall rate  
 491 from top-5 to top-30 than the other three models,

Models	Hits@1	F1
LLaMA2-7B-Chat	33.6	13.5
Ours	58.6	20.1
Ours (continual training)	75.4	50.1

Table 3: Results on MetaQA-3hop.

492 significantly surpassing the original BGE model.  
 493 Comparing the performance of our model without  
 494 CoT and without subgraph information, we find  
 495 that subgraph information is more crucial for the  
 496 retrieval model, consistent with the results of the  
 497 ablation study in Section 5.1.

## 498 5.3 The Efficiency of YAML Format KG

499 As analyzed in Section 3.3, adopting the YAML  
 500 format with simple syntax to represent KGs instead  
 501 of the traditional triplet format can reduce token  
 502 redundancy. To quantitatively assess how much  
 503 redundancy YAML can eliminate, we have calcu-  
 504 lated the average number of KG tokens required  
 505 per question by selecting knowledge graphs con-  
 506 structed from knowledge retrieved by our search  
 507 engine on both WebQSP and CWQ datasets. For  
 508 WebQSP, using triples to represent the KG requires  
 509 an average of 532.6 tokens per question; if we  
 510 use the YAML format, the average token drops to  
 511 384.2, thus reducing token redundancy by nearly  
 512 28%. For CWQ, replacing triples with YAML re-  
 513 duces the average token count of KGs from 534.3 to  
 514 401.4, a compression of nearly 25%. In a scenario  
 515 where budget resources are constrained, minimiz-  
 516 ing the representation of tokens in a knowledge  
 517 graph by using YAML allows those resources to be  
 518 repurposed towards combining additional examples  
 519 or recalling more retrieved information, aiming to  
 520 achieve further performance enhancements.

## 521 5.4 Transferring to Other KGs

522 To further validate the transferability of our method  
 523 to other KGs, we choose the MetaQA-3hop dataset  
 524 (Zhang et al., 2018), which is based on the Wiki-  
 525 data KG. We continue training our method on mod-  
 526 els trained on the Freebase KG and the WebQSP,  
 527 CWQ datasets. We construct 35k samples from  
 528 the Wikidata KG for KG instruction tuning, sam-  
 529 ple 75k samples from the MetaQA-3hop training  
 530 set for training the retrieval module, and use 60k  
 531 reasoning processes as the training data for KG re-  
 532asoning. Training details are consistent with those  
 533 described in Section 4.1. As Table 3 shows, the  
 534 original LLaMA2-7b-Chat performs poorly on the

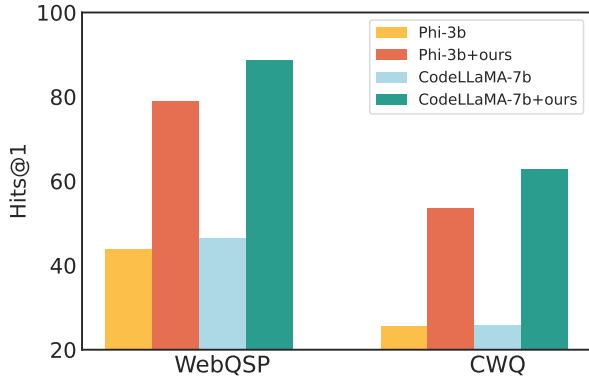


Figure 4: Experimental results on Phi2 and CodeLLaMA models.

dataset. Even without further training on Wikidata and MetaQA, our model achieves a 25% improvement in Hits@1. After continual training, the performance of our method continues to rise, reaching 75% for Hits@1. It indicates that our method is equipped with KG retrieval, comprehension and reasoning capabilities independent of specific datasets, and has the huge potential for transferability.

### 5.5 Applying to Other Models

To verify the generalizability of our proposed method, we apply our method on two other different models, CodeLLaMA-7b-Instruct<sup>6</sup> (Rozière et al., 2024) and Phi2-3b<sup>7</sup> (Li et al., 2023b). As shown in Figure 4, our method has significantly improved the performance of these two models on the KGQA task. For Phi2 and CodeLLaMA, our method has achieved an average improvement of 30% and 40% on the two datasets, respectively. Although CodeLLaMA is slightly inferior to LLaMA2-7b-Chat, it still achieves performance comparable to RoG. Phi2, with only half the number of parameters compared to the other two models, lags significantly behind in performance, also reaching the level of UniKGQA and ChatGPT.

We observe that the performance differences among the original three models on KGQA tasks are not significant. This phenomenon offers new insights for selecting a foundational model for KGQA in practice: firstly, within resource limits, choose models with larger parameters to fully learn and utilize KG capabilities; secondly, choose models with stronger reasoning abilities.

<sup>6</sup><https://huggingface.co/codellama/CodeLlama-7b-Instruct-hf>

<sup>7</sup><https://huggingface.co/microsoft/phi-2>

### 5.6 Error Analysis and Case Study

We conduct error analysis to further explore the strengths and weaknesses of our proposed method in CWQ. Firstly, we categorize all the model’s predictions into three levels based on evaluation metrics: perfect predictions (Hit@1=1 and F1=1), imperfect predictions (Hit@1=1 and 0<F1<1), and completely wrong predictions (Hit@1=0 and F1=0). In our method to the CWQ dataset, the predictions for these three levels are distributed as 39.1%, 29.6%, and 31.3%, respectively. Furthermore, we observe that the ground truths for perfect predictions are short, with 86% containing only a single answer. This indicates that our method still struggles to achieve perfect predictions for complex questions with multiple answers. For queries where the model produces imperfect predictions, their recall (80.7%) is significantly higher than their precision (50.6%), indicating that our model still suffers from hallucination. While it can provide correct answers, it may also be misled by irrelevant information retrieved, resulting in inaccurate answers. Finally, we find that queries with completely wrong predictions also have fewer answers, but these queries exhibited inferior retrieval quality compared to those with perfect predictions. We consider the knowledge contained in the final step of every path leading to an answer to be the most crucial; queries with perfect predictions attained a recall rate for this knowledge of 81.1%, whereas those with entirely incorrect predictions only reached 42.4%. Thus, enhancing the retrieval model’s ability to handle complex queries and improving the reasoning model’s resistance to irrelevant retrieved content are promising directions for further advancing the performance of LLM in KGQA tasks. We show cases in Appendix C.

## 6 Conclusion

In this paper, we propose a method combining explainable knowledge graphs with large language models to enhance complex reasoning capabilities. Our method includes a KG retrieval model and a KG reasoning model. We integrate reasoning processes and subgraph information for better KG retrieval. We employ a novel KG representation and KG-related tuning for the reasoning model to learn to understand and reason with KG. Experimental results on two challenging KGQA tasks show that our method outperforms existing strong baselines and the SOTA model.



## 618 Limitations

619 Although our proposed method has made signifi-  
620 cant progress in KGQA, there are still some limita-  
621 tions:

- 622 • Due to computational resource constraints, we  
623 conduct experiments only on LLMs below 10B  
624 parameters, lacking investigation into larger mod-  
625 els (such as LLaMA2-13B and 70B), other archi-  
626 tectures (such as RWKV and Mixtral families).
- 627 • Our method fine-tunes LLMs with full-parameter,  
628 which is impractical in many low-resource set-  
629 tings. In future work, we plan to utilize efficient  
630 fine-tuning techniques such as LoRA, and com-  
631 pare its effectiveness with the current results.
- 632 • We validate the efficacy of our method only on  
633 two KGQA tasks. To more convincingly demon-  
634 strate that our approach enables LLMs to lever-  
635 age KG for reasoning, we will incorporate addi-  
636 tional tasks and datasets in our future work.

## 637 References

638 Jinheon Baek, Alham Fikri Aji, Jens Lehmann, and  
639 Sung Ju Hwang. 2023. [Direct fact retrieval from  
640 knowledge graphs without entity linking](#). In *Proceed-  
641 ings of the 61st Annual Meeting of the Association for  
642 Computational Linguistics (Volume 1: Long Papers)*,  
643 pages 10038–10055, Toronto, Canada. Association  
644 for Computational Linguistics.

645 Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim  
646 Sturge, and Jamie Taylor. 2008. [Freebase: a col-  
647 laboratively created graph database for structuring  
648 human knowledge](#). In *Proceedings of the 2008 ACM  
649 SIGMOD International Conference on Management  
650 of Data*, SIGMOD '08, page 1247–1250, New York,  
651 NY, USA. Association for Computing Machinery.

652 Sébastien Bubeck, Varun Chandrasekaran, Ronen El-  
653 dan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Pe-  
654 ter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg,  
655 Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro,  
656 and Yi Zhang. 2023. [Sparks of artificial general in-  
657 telligence: Early experiments with gpt-4](#).

658 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,  
659 Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo,  
660 Meng Wang, and Haofen Wang. 2024. [Retrieval-  
661 augmented generation for large language models: A  
662 survey](#).

663 Claire Gardent, Anastasia Shimorina, Shashi Narayan,  
664 and Laura Perez-Beltrachini. 2017. [Creating training  
665 corpora for NLG micro-planners](#). In *Proceedings  
666 of the 55th Annual Meeting of the Association for  
667 Computational Linguistics, ACL 2017, Vancouver,*

*Canada, July 30 - August 4, Volume 1: Long Pa-  
pers*, pages 179–188. Association for Computational  
Linguistics.

668 Yu Gu and Yu Su. 2022. [ArcaneQA: Dynamic program  
669 induction and contextualized encoding for knowl-  
670 edge base question answering](#). In *Proceedings of  
671 the 29th International Conference on Computational  
672 Linguistics*, pages 1718–1731, Gyeongju, Republic  
673 of Korea. International Committee on Computational  
674 Linguistics.

675 Gaole He, Yunshi Lan, Jing Jiang, Wayne Xin Zhao, and  
676 Ji-Rong Wen. 2021. [Improving multi-hop knowledge  
677 base question answering by learning intermediate  
678 supervision signals](#). In *Proceedings of the 14th ACM  
679 International Conference on Web Search and Data  
680 Mining, WSDM '21*, page 553–561, New York, NY,  
681 USA. Association for Computing Machinery.

682 Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla,  
683 Thomas Laurent, Yann LeCun, Xavier Bresson, and  
684 Bryan Hooi. 2024. [G-retriever: Retrieval-augmented  
685 generation for textual graph understanding and ques-  
686 tion answering](#).

687 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,  
688 Zhangyin Feng, Haotian Wang, Qianglong Chen,  
689 Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting  
690 Liu. 2023. [A survey on hallucination in large lan-  
691 guage models: Principles, taxonomy, challenges, and  
692 open questions](#).

693 Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin  
694 Zhao, and Ji-Rong Wen. 2023a. [StructGPT: A gen-  
695 eral framework for large language model to reason  
696 over structured data](#). In *Proceedings of the 2023 Con-  
697 ference on Empirical Methods in Natural Language  
698 Processing*, pages 9237–9251, Singapore. Associa-  
699 tion for Computational Linguistics.

700 Jinhao Jiang, Kun Zhou, Xin Zhao, and Ji-Rong Wen.  
701 2023b. [UniKGQA: Unified retrieval and reasoning  
702 for solving multi-hop question answering over knowl-  
703 edge graph](#). In *The Eleventh International Confer-  
704 ence on Learning Representations*.

705 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao  
706 Fu, Kyle Richardson, Peter Clark, and Ashish Sab-  
707 harwal. 2023. [Decomposed prompting: A modular  
708 approach for solving complex tasks](#). In *The Eleventh  
709 International Conference on Learning Representa-  
710 tions*.

711 Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi.  
712 2023. [KG-GPT: A general framework for reasoning  
713 on knowledge graphs using large language models](#).  
714 In *Findings of the Association for Computational Lin-  
715 guistics: EMNLP 2023*, pages 9410–9421, Singapore.  
716 Association for Computational Linguistics.

717 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying  
718 Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
719 Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Effi-  
720 cient memory management for large language model  
721 serving with pagedattention](#). In *Proceedings of the*

725	<i>ACM SIGOPS 29th Symposium on Operating Systems Principles.</i>	
726		
727	Yunshi Lan and Jing Jiang. 2020. <a href="#">Query graph generation for answering multi-hop complex questions from knowledge bases.</a> In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 969–974, Online. Association for Computational Linguistics.	
728		
729		
730		
731		
732		
733	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. <a href="#">Retrieval-augmented generation for knowledge-intensive nlp tasks.</a> In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 9459–9474. Curran Associates, Inc.	
734		
735		
736		
737		
738		
739		
740		
741	Tianle Li, Xueguang Ma, Alex Zhuang, Yu Gu, Yu Su, and Wenhui Chen. 2023a. <a href="#">Few-shot in-context learning on knowledge base question answering.</a> In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6966–6980, Toronto, Canada. Association for Computational Linguistics.	
742		
743		
744		
745		
746		
747		
748	Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023b. <a href="#">Textbooks are all you need ii: phi-1.5 technical report.</a>	
749		
750		
751		
752	Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2020. <a href="#">K-bert: Enabling language representation with knowledge graph.</a> <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(03):2901–2908.	
753		
754		
755		
756		
757	Linho Luo, Yuan-Fang Li, Reza Haf, and Shirui Pan. 2024. <a href="#">Reasoning on graphs: Faithful and interpretable large language model reasoning.</a> In <i>The Twelfth International Conference on Learning Representations</i> .	
758		
759		
760		
761		
762	Alexander Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. 2016. <a href="#">Key-value memory networks for directly reading documents.</a> In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1400–1409, Austin, Texas. Association for Computational Linguistics.	
763		
764		
765		
766		
767		
768		
769	OpenAI. 2022. Chatgpt: Optimizing language models for dialogue. <i>Open AI, blog</i> .	
770		
771	OpenAI. 2023. Gpt-4 technical report. <i>ArXiv</i> , abs/2303.08774.	
772		
773	Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jipapu Wang, and Xindong Wu. 2024. <a href="#">Unifying large language models and knowledge graphs: A roadmap.</a> <i>IEEE Transactions on Knowledge and Data Engineering</i> , page 1–20.	
774		
775		
776		
777		
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. <a href="#">Code llama: Open foundation models for code.</a>	778 779 780 781 782 783 784 785 786 787
	Jiaxin Shi, Shulin Cao, Lei Hou, Juanzi Li, and Hanwang Zhang. 2021. <a href="#">TransferNet: An effective and transparent framework for multi-hop question answering over relation graph.</a> In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4149–4158, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	788 789 790 791 792 793 794 795
	Haitian Sun, Tania Bedrax-Weiss, and William Cohen. 2019. <a href="#">PullNet: Open domain question answering with iterative retrieval on knowledge bases and text.</a> In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2380–2390, Hong Kong, China. Association for Computational Linguistics.	796 797 798 799 800 801 802 803 804
	Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William Cohen. 2018. <a href="#">Open domain question answering using early fusion of knowledge bases and text.</a> In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4231–4242, Brussels, Belgium. Association for Computational Linguistics.	805 806 807 808 809 810 811 812
	Yawei Sun, Lingling Zhang, Gong Cheng, and Yuzhong Qu. 2020. <a href="#">Sparqa: Skeleton-based semantic parsing for complex questions over knowledge bases.</a> <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(05):8952–8959.	813 814 815 816 817
	Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, Yanbin Zhao, Yuxiang Lu, Weixin Liu, Zhihua Wu, Weibao Gong, Jianzhong Liang, Zhizhou Shang, Peng Sun, Wei Liu, Xuan Ouyang, Dianhai Yu, Hao Tian, Hua Wu, and Haifeng Wang. 2021. <a href="#">Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation.</a>	818 819 820 821 822 823 824 825
	Alon Talmor and Jonathan Berant. 2018. <a href="#">The web as a knowledge-base for answering complex questions.</a> In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 641–651, New Orleans, Louisiana. Association for Computational Linguistics.	826 827 828 829 830 831 832 833
	Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang	834 835

836	Xiong. 2023. <a href="#">Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering.</a>	<i>Linguistics (Volume 1: Long Papers)</i> , pages 5773–5784, Dublin, Ireland. Association for Computational Linguistics.	893
837			894
838			895
839	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. <a href="#">Chain-of-thought prompting elicits reasoning in large language models.</a> In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 24824–24837. Curran Associates, Inc.	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. <a href="#">Siren’s song in the ai ocean: A survey on hallucination in large language models.</a>	896
840			897
841			898
842			899
843			900
844			901
845			
846	Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. <a href="#">The dawn of llms: Preliminary explorations with gpt-4v(ision).</a>	Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander J Smola, and Le Song. 2018. Variational reasoning for question answering with knowledge graph. In <i>AAAI</i> .	902
847			903
848			904
849			905
850	Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. <a href="#">QA-GNN: Reasoning with language models and knowledge graphs for question answering.</a> In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 535–546, Online. Association for Computational Linguistics.	Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. <a href="#">ERNIE: Enhanced language representation with informative entities.</a> In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1441–1451, Florence, Italy. Association for Computational Linguistics.	906
851			907
852			908
853			909
854			910
855			911
856			912
857			
858	Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. 2024. <a href="#">Language is all a graph needs.</a> In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 1955–1973, St. Julian’s, Malta. Association for Computational Linguistics.		
859			
860			
861			
862			
863			
864	Xi Ye, Semih Yavuz, Kazuma Hashimoto, Yingbo Zhou, and Caiming Xiong. 2022. <a href="#">RNG-KBQA: Generation augmented iterative ranking for knowledge base question answering.</a> In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6032–6043, Dublin, Ireland. Association for Computational Linguistics.		
865			
866			
867			
868			
869			
870			
871			
872	Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. 2015. <a href="#">Semantic parsing via staged query graph generation: Question answering with knowledge base.</a> In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1321–1331, Beijing, China. Association for Computational Linguistics.		
873			
874			
875			
876			
877			
878			
879			
880			
881	Donghan Yu, Sheng Zhang, Patrick Ng, Henghui Zhu, Alexander Hanbo Li, Jun Wang, Yiqun Hu, William Yang Wang, Zhiguo Wang, and Bing Xiang. 2023. <a href="#">DecAF: Joint decoding of answers and logical forms for question answering over knowledge bases.</a> In <i>The Eleventh International Conference on Learning Representations</i> .		
882			
883			
884			
885			
886			
887			
888	Jing Zhang, Xiaokang Zhang, Jifan Yu, Jian Tang, Jie Tang, Cuiping Li, and Hong Chen. 2022. <a href="#">Subgraph retrieval enhanced model for multi-hop knowledge base question answering.</a> In <i>Proceedings of the 60th Annual Meeting of the Association for Computational</i>		
889			
890			
891			
892			

## A Datasets

Table 4 show the characteristics of WebQSP and CWQ. To narrow the scope of retrieval, we follow previous work (Luo et al., 2024; Jiang et al., 2023b) and construct a subgraph of Freebase by extracting all triples within the max reasoning hops of question entities in WebQSP and CWQ.

Dataset	WebQSP	CWQ
#Train	2,848	27,639
#Valid	250	3,519
#Test	1,639	3,519
#Max hop	2	5

Table 4: Characteristics of datasets

## B KG-oriented Instruction

We list all instruction templates used in KG-oriented instruction tuning in Table 5.

## C Case Study

We provide the generation cases from the original LLaMA2-7b-chat and our model below, where we use red to indicate the factual error. We can observe that our model reduced factual errors in the reasoning process by utilizing knowledge graph knowledge retrieved. However, it may still provide irrelevant answers due to the interference of unrelated retrieval information.

### Case 1

#### Question:

What is the Title of the national anthem for the country where the Vyakta horse breed originated in?

**Ground truth:** ['National Anthem of Russia']

#### LLaMA2-7b-chat:

1. Identify the key components of the question. The question asks about the national anthem of a country and the origin of a horse breed.
2. Determine the country associated with the Vyakta horse breed: **The Vyakta horse breed originated in India.**
3. Find out the national anthem of India: The national anthem of India is "Jana Gana Mana."  
Output: ['Jana Gana Mana']

#### Ours:

- Step 1: Identify the Country of Origin. The Vyakta horse breed originated in Russia.
- Step 2: Find the National Anthem. The national anthem of Russia is "God Save the Tsar".  
Output: ['National Anthem of Russia']

### Case 2

#### Question:

What currency is used in the country with Nobel Patria, tu hermosa as its national anthem?

**Ground truth:** ['Costa Rican colón']

#### LLaMA2-7b-chat:

1. Identify the components of the question: The question asks for the currency used in a country with a specific national anthem.
2. Determine the country: **The country with the national anthem "Nobel Patria, tu hermosa" is Venezuela.**
3. Find out the currency used in Venezuela: The currency used in Venezuela is the Bolivar (VEF).  
Output: ['Bolívar (VEF)']

#### Ours:

- Step 1: Identify the National Anthem. The national anthem in question is "Noble patria, tu hermosa bandera".
- Step 2: Identify the Country. The country that uses the national anthem "Noble patria, tu hermosa bandera" is Costa Rica.
- Step 3: Identify the Currency. The currency used in Costa Rica is the Costa Rican colón.  
Output: ['Costa Rican colón']

### Case 3

#### Question:

Which sports are popular in Spain, and what is the name of the team that plays in Donostia/San Sebastian?

**Ground truth:** ['San Sebastián Gipuzkoa BC']

#### LLaMA2-7b-chat:

1. Identify the country: Spain.
2. Determine the popular sports in Spain: **Football (soccer) is the most popular sport in Spain.**
3. Find the name of the team that plays in Donostia/San Sebastian: The team that plays in Donostia/San Sebastian is Real Sociedad.  
Output: ['Football', 'Real Sociedad']

#### Ours:

- Step 1: Identify the Location. The location of interest is Donostia/San Sebastián.
- Step 2: Find Teams in the Location. The teams located in Donostia/San Sebastián are Real Sociedad and San Sebastián Gipuzkoa BC.
- Step 3: Identify Sports. The sports teams in the location are Real Sociedad (associated with football) and San Sebastián Gipuzkoa BC (associated with basketball).
- Step 4: Find Other Teams. The teams associated with Spain are the Spain national football team, CB Estudiantes, Basket Zaragoza 2002, Saski Baskonia, Spain men's national volleyball team, Spain women's national volleyball team, CB 1939 Canarias, and Baloncesto Málaga.
- Step 5: Identify Popular Sports. **The popular sports in Spain are football, basketball, and volleyball.**  
Output: ['Real Sociedad', 'San Sebastián Gipuzkoa BC']

Task	ID	Instruction
Entity	1	Please predict the entity represented by <mask> based on the one-hop relationships in the knowledge graph. Input: {Input} Output: {Output}
	2	Based on the one-hop relationships in the knowledge graph, infer the entity represented by <mask>. Input: {Input} Output: {Output}
	3	Make a prediction about the masked entity, using the one-hop relationships in the knowledge graph as a reference. Input: {Input} Output: {Output}
Relationship	1	Please recognize the relationship between the two entities. Knowledge Graph: {KG} Input: {Input} Output: {Output}
	2	Please predict the relationship between the two entities. There are some one-hop information of these entities: {KG} Input: {Input} Output: {Output}
	3	Make a prediction about the relationship, using the one-hop relationships in the knowledge graph as a reference. {KG} Input: {Input} Output: {Output}
Graph2text	1	Please deeply understand the following knowledge graph, and then convert them into a coherent sentence. Input: {Input} Output: {Output}
	2	Given these knowledge graph, please deeply write a paragraph that integrates the information contained in them. Input: {Input} Output: {Output}
	3	Compose an informative report using the information from these knowledge graph. Input: {Input} Output: {Output}
Text2graph	1	Please extract all entities and relationships in the sentence. Input: {Input} Output: {Output}
	2	Given the sentence, please extract a knowledge graph that integrates the information contained in them. Input: {Input} Output: {Output}
	3	Please deeply understand the following sentence, and then generate a knowledge graph. Input: {Input} Output: {Output}

Table 5: Instructions of the KG-related tasks.