

# Retrieve-Rewrite-Answer: A KG-to-Text Enhanced LLMs Framework for Knowledge Graph Question Answering

Yike Wu\*  
yike.wu@seu.edu.cn  
Southeast University  
Nanjing, Jiangsu, China

Nan Hu\*  
nanhu@seu.edu.cn  
Southeast University  
Nanjing, Jiangsu, China

Sheng Bi  
bisheng@seu.edu.cn  
Southeast University  
Nanjing, Jiangsu, China

Guilin Qi  
gqi@seu.edu.cn  
Southeast University  
Nanjing, Jiangsu, China

Jie Ren  
renjie@zhejianglab.com  
Zhejiang Lab  
Hangzhou, Zhejiang, China

Anhuan Xie  
xieanhuan@zhejianglab.com  
Zhejiang Lab  
Hangzhou, Zhejiang, China

Wei Song<sup>†</sup>  
weisong@zhejianglab.com  
Zhejiang Lab  
Hangzhou, Zhejiang, China

LLM이 여전히 KIT와 long-text query에 대한  
정답을 제대로 생성하기 못함.

## ABSTRACT

Despite their competitive performance on knowledge-intensive tasks, large language models (LLMs) still have limitations in memorizing all world knowledge especially long tail knowledge. In this paper, we study the KG-augmented language model approach for solving the knowledge graph question answering (KGQA) task that requires rich world knowledge. Existing work has shown that retrieving KG knowledge to enhance LLMs prompting can significantly improve LLMs performance in KGQA. However, their approaches lack a well-formed verbalization of KG knowledge, i.e., they ignore the gap between KG representations and textual representations. To this end, we propose an answer-sensitive KG-to-Text approach that can transform KG knowledge into well-textualized statements most informative for KGQA. Based on this approach, we propose a KG-to-Text enhanced LLMs framework for solving the KGQA task. Experiments on several KGQA benchmarks show that the proposed KG-to-Text augmented LLMs approach outperforms previous KG-augmented LLMs approaches regarding answer accuracy and usefulness of knowledge statements.<sup>1</sup>

## CCS CONCEPTS

• Computing methodologies → Artificial intelligence; Natural language processing; Natural language generation.

## KEYWORDS

Knowledge Graph Augmented LLMs, Knowledge Graph Question Answering, KG-to-Text

## 1 INTRODUCTION

Large language models (LLMs) are becoming increasingly popular in natural language processing for their superior competence in various applications. Although LLMs demonstrate remarkable capabilities in zero-shot scenarios, their performances on several

knowledge-intensive tasks are insufficiently satisfactory [20]. This reveals the enormous parameters in LLMs cannot store all the world's knowledge. Several researches indicate that LLMs still suffer from issues like hallucinations and factual inaccuracy when answering questions [12, 26]. Specifically, LLMs perform inadequately in the knowledge-intensive task KGQA [9, 34].

To solve this problem, recent work attempts to enhance LLMs with external knowledge. A line of work [10, 19, 38, 41] involves continual pre-training LLMs on extensive corpora. Nevertheless, this method requires a significant amount of textual data, computational resources, and time investment. Some previous work has attempted to explicitly enhance the LLMs' performance on downstream tasks by leveraging external knowledge, such as knowledge graphs and web contents [11, 14, 18]. This approach is employed to address the model's deficiency in factual knowledge. Inspired by this, other work [1, 32] constructs knowledge-augmented prompt by prepending question-related factual information to the question to enrich the knowledge of LLMs in a more direct fashion. Although this approach proves to be successful and cost-effective, it ignores the importance of knowledge representation.

In this paper, we summarize two knowledge representation formats employed in previous work: triple-form text and free-form text. As shown in Figure 1, triple-form text involves a simple linear concatenation of triples, transforming them into structured text. Free-form text converts triples into semantically coherent textual descriptions based on the KG-to-Text method. Further, we propose a KG-to-Text enhanced framework, Retrieve-Rewrite-Answer to strengthen the performance of LLMs on KGQA. As shown in Figure 2, compared with previous work that answers questions in a Retrieve-then-Answer fashion, our framework adopts a Rewriter module to transform the retrieved triples into textual descriptions. The core of this framework lies in the task-driven KG-to-Text method. Our designed method is answer-sensitive and can transform question-related triples into textual knowledge which is most informative to KGQA. Compared with previous work that simply adopts an off-the-shelf KG-to-Text model, we fine-tune open-source

P1. Hallucination

P2.

P3.

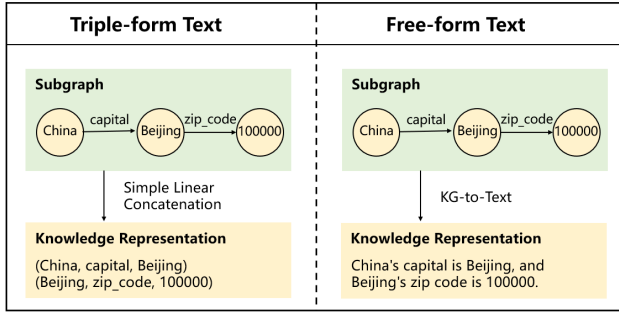
\*Both authors contributed equally to this research.

<sup>†</sup>Corresponding author

<sup>1</sup>Our code is available at <https://github.com/wuyike2000/Retrieve-Rewrite-Answer>

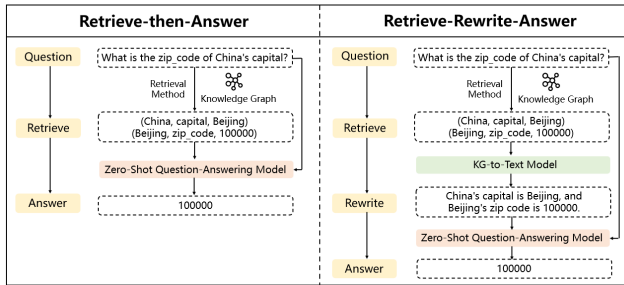
off-the-shelf = 기성품의, 미리 만들어 놓은,  
특히 표준화된 방법으로 만들어 둔

1. Triple 형식 텍스트는 트리플을 단순히 선형적으로 연결하여 구조화된 텍스트로 변환함.
2. 자유 형식 텍스트는 KG-to-Text를 기반으로 트리플을 의미론적으로 일관된 텍스트 설명으로 변환함.
3. RRW => KGQA에 대한 LLM의 성능을 강화함.



**Figure 1: Two forms of knowledge representation: triple-form text and free-form text. Triple-form text transforms subgraph to text via simple linear concatenation while free-form text adopts the KG-to-Text method to generate a semantically coherent textual description.**

LLMs on the KG-to-Text corpus to generate knowledge descriptions beneficial to KGQA. However, the major challenge is the lack of KG-to-Text annotated data in existing KGQA benchmarks. Therefore, we propose an automatic corpus generation method for generating high-quality graph-text pairs based on the feedback of LLMs.



**Figure 2: Comparison between Retrieve-then-Answer framework and our proposed Retrieve-Rewrite-Answer framework. We use a rewriter module to transform the retrieved subgraphs into textual descriptions.**

The contributions of this paper are summarized as follows:

- We propose **Retrieve-Rewrite-Answer**, a KG-to-Text enhanced LLMs framework for KGQA. In comparison to previous KG-augmented LLMs frameworks, the most significant innovation in our framework lies in the **rewrite module**. This module uses **fine-tuned LLMs as KG-to-Text models** to convert retrieved subgraphs into well-textualized statements most informative for KGQA.
- To address the issue of scarcity of KG-to-Text corpus annotations, we devise an **automatic KG-to-Text corpus generation** method. We extract **question-related subgraphs** and utilize **ChatGPT as a tool for corpus generation**. Based on the feedback from question-answering LLMs, we generate answer-sensitive knowledge descriptions for the construction of KG-to-Text labeled data. We **fine-tune several open-source**

LLMs on the generated corpus and investigate the impact of textual knowledge generated by different LLMs on KGQA.

- We evaluate our framework on four KGQA benchmarks. Experimental results show that our framework outperforms previous KG-augmented methods across several LLMs, which demonstrates its effectiveness. Besides, we investigate the influence of different knowledge representation formats on KGQA and prove that the knowledge produced by our framework is most beneficial.

## 2 PRELIMINARIES

**Knowledge Graph (KG)** is a collection of triples  $(s, r, o)$  composed of subject  $s$ , relation  $r$  and object  $o$ , denoted by  $G = \{(s, r, o) | s, o \in E, r \in R\}$ , where  $E$  and  $R$  denote the entity set and relation set.

**KG-to-Text** is a natural language generation technique based on KG. Given a subgraph  $G' = \{(s, r, o) | s, o \in E, r \in R\}$  from KG  $G$ , the objective of KG-to-Text is to generate a text sequence  $X = (x_1, x_2, \dots, x_n)$  that is **semantically coherent** with subgraph  $G'$ .

**Knowledge Graph Question Answering (KGQA)** is the task of answering natural language questions based on a set of facts over KGs. Given a question  $q$  and a topic entity  $e_h$ , the task is to generate an answer  $a$  that can correctly respond to the question.

## 3 RELATED WORK

### 3.1 KG-Augmented LLM for KGQA

Despite LLMs' pre-training on massive corpora, they still suffer from hallucinations, factual inaccuracy, and outdated knowledge when it comes to KGQA. Recent efforts have aimed to harness KGs to enhance the capabilities of LLMs on KGQA [1, 32]. In these studies, question-related triples are extracted from KG and transformed into textual format using techniques such as linear verbalization or off-the-shelf KG-to-Text models. **The textual representation of triples and the question are transformed into the knowledge-augmented prompt via predefined templates. The prompt is then fed into the question-answering LLM to produce more reliable answers.**

While these studies demonstrate the effectiveness of this strategy, they **ignore the impact of the knowledge representation format on the performance of LLMs on KGQA**. In this study, we devise a KG-to-Text corpus generation method to produce high-quality annotations. Then we utilize LLMs fine-tuned on the corpus to transform question-related subgraphs into knowledge text, which can further elevate the performance of LLMs on KGQA.

### 3.2 KG-to-Text

Recent studies in KG-to-Text can be classified into two main approaches: 1) based on graph neural networks (GNNs) [40]: In order to **preserve the structural information** of subgraphs during encoding, researchers focus on designing more complex encoders to obtain enhanced subgraph representation. Some work has developed graph-structured encoders based on GNNs [6, 21] due to the robust capabilities of GNNs in encoding graph-structured data. However, **GNNs are restricted by local processing nature, as they iteratively compute node representations based on neighboring nodes' features**. To overcome this limitation, alternative efforts have employed the Transformer-based architecture [39] in designing encoders [3, 15, 46]. This approach allows for the establishment

Suggest

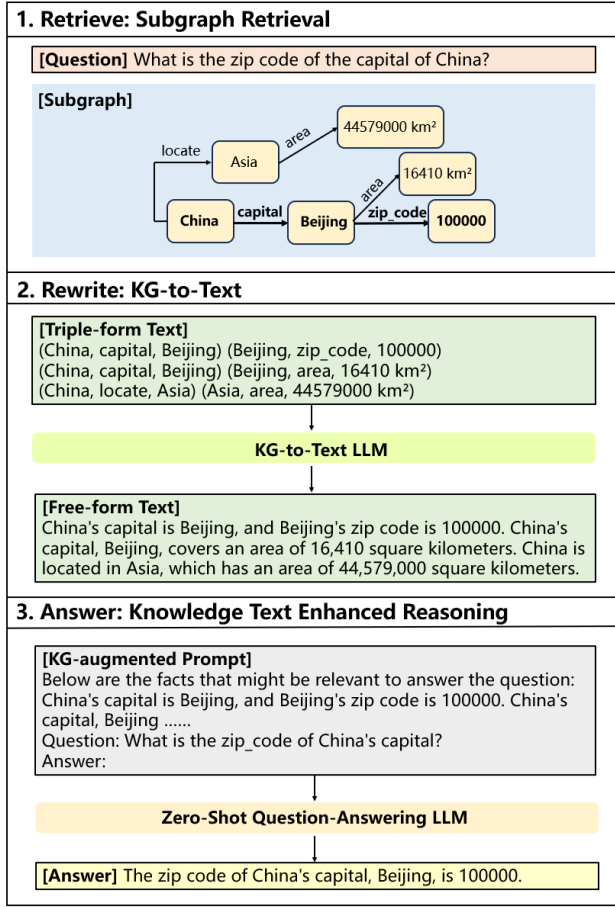


Figure 3: Our proposed framework, Retrieve-Rewrite-Answer has three steps: subgraph retrieval, KG-to-Text, and knowledge text enhanced reasoning.

of connections between any two nodes within the graph, thereby addressing the limitations of GNNs. 2) based on pre-trained language models (PLMs): With the development of large-scale generative PLMs such as BART [17], T5 [24] and GPT [23], recent work applies these models to KG-to-Text, modeling it as an end-to-end generation task [7, 13, 25]. These studies involve modifications to the model architecture and the introduction of pre-training tasks to enhance the ability to extract structural information. In this work, we follow the second approach and fine-tune open-source LLMs on the KG-to-Text corpus.

## 4 METHODS

### 4.1 Retrieve-Rewrite-Answer

As shown in Figure 3, our proposed Retrieve-Rewrite-Answer framework contains three steps: subgraph retrieval, KG-to-Text, and knowledge text enhanced reasoning.

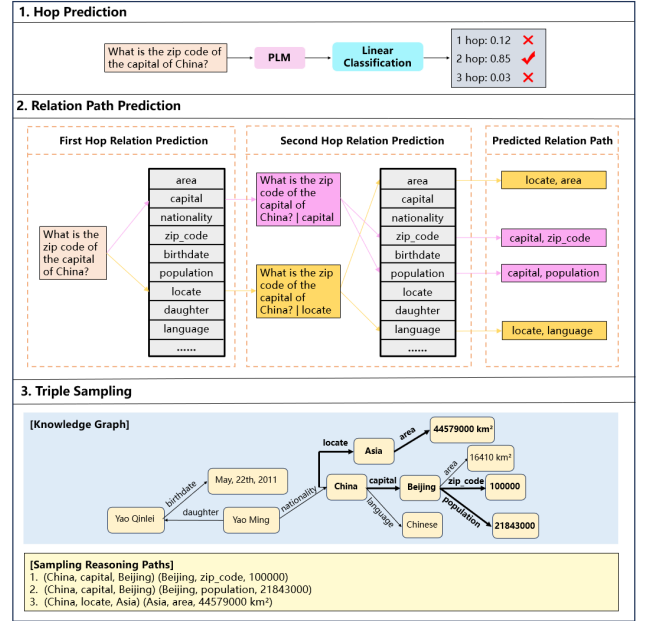


Figure 4: Our retrieve module consists of three steps: hop prediction, relation path prediction, and triple sampling.

4.1.1 *Retrieve: Subgraph Retrieval.* Our retrieve module consists of three steps: hop prediction, relation path prediction, and triple sampling, as illustrated in Figure 4.

**Hop Prediction.** The aim of this step is to predict the hop number of the question, which is used to predict the relation path in the next step. We model the hop prediction as a classification task based on PLMs. Given the question  $q$ , we employ PLMs to encode the question  $q$  and obtain the vector representation  $q_v$ :

$$q_v = PLM(q) \quad (1)$$

The representation  $q_v$  is then fed into the linear classification layer to predict the probability distribution  $D'_h$  of potential hop numbers  $h_1, h_2, \dots, h_H$ :

$$D'_h = [d'_{h_1}, d'_{h_2}, \dots, d'_{h_H}] = Linear(q_v) \quad (2)$$

where  $d'_{h_c}$  is the probability of the hop number  $h_c$  given the question representation  $q_v$ :

$$d'_{h_c} = P(h_c | q_v), \quad c = 1, 2, \dots, H \quad (3)$$

The hop number  $h$  with the highest probability is selected as the predicted result.

$$h = \arg \max_{h_c} d'_{h_c}, \quad c = 1, 2, \dots, H \quad (4)$$

During training, ground truth distribution  $D_h$  is represented as a one-hot vector with a probability of 1 for the true hop number  $h_{gold}$  and probabilities of 0 for other hop numbers:

$$D_h = [d_{h_1}, d_{h_2}, \dots, d_{h_H}] \quad (5)$$

$$d_{h_c} = \begin{cases} 1, & h_c = h_{gold} \\ 0, & h_c \neq h_{gold} \end{cases} \quad c = 1, 2, \dots, H \quad (6)$$

The predicted distribution  $D'_h$  is penalized for being different from the ground truth distribution  $D_h$  using **cross entropy loss**  $L_{CE}$ :

$$L_{CE} = -D_h \log D'_h = - \sum_{c=1}^H d_{h_c} \log d'_{h_c} \quad (7)$$

Predicted distribution  $\leftarrow$   $D'_h$   
Ground Truth  $\rightarrow$   $D_h$

$L_{CE}$  is utilized to update the model's parameters.

**Relation Path Prediction.** Given the question  $q$  and the predicted hop number  $h$ , we perform  $h$ -step prediction, with each step corresponding to one hop relation. At step  $t$ , we predict the  $t$ -th hop relation based on the predicted  $(t-1)$  hop relation paths and the question  $q$  via PLMs as a classification task. Specifically, for each predicted relation path, we sample  $K$  candidate relations for the next step. In step 1, we encode the question  $q$  into a vector representation  $q_v$ . This vector is then passed through a linear classification layer to compute the distribution  $D'_{r,1}$  of  $R$  relations in KG:

$$D'_{r,1} = [d'_{r_1}, d'_{r_2}, \dots, d'_{r_R}] = \text{Linear}(q_v) \quad (8)$$

given question에 따른 R번째 relation의 확률

where  $d'_c$  is the probability of relation  $r_c$  given the question representation  $q_v$ :

$$d'_c = P(r_c | q_v), \quad c = 1, 2, \dots, R \quad (9)$$

question의 임베딩

We select top  $K$  relations with the highest probabilities as one hop relation paths  $p_1$ .

In the following step  $t$  ( $t > 1$ ), the relation path  $p_{t-1,i}$  in the  $(t-1)$  hop relation paths  $p_{t-1}$  can be represented as:

$$p_{t-1,i} = r_{i,1} | r_{i,2} | \dots | r_{i,t-1}, \quad i = 1, 2, \dots, K^{t-1} \quad (10)$$

We concatenate the question  $q$  and the relation path  $p_{t-1,i}$  with “|” as the input sequence  $Q_t$ :

$$Q_t = q | r_{i,1} | r_{i,2} | \dots | r_{i,t-1} \quad (11)$$

$Q_t$  is encoded into vector representation  $Q_{t,v}$  via PLMs and passes through a linear classification layer to calculate the relation distribution  $D'_{r,t}$  across  $R$  relations in KG:

$$Q_{t,v} = \text{PLM}(Q_t) \quad (12)$$

$$D'_{r,t} = [d'_{r_1}, d'_{r_2}, \dots, d'_{r_R}] = \text{Linear}(Q_{t,v}) \quad (13)$$

where  $d'_c$  is the probability of relation  $r_c$  given the input sequence representation  $Q_{t,v}$ :

$$d'_c = P(r_c | Q_{t,v}), \quad c = 1, 2, \dots, R \quad (14)$$

We retain top  $K$  relations with the highest probabilities as the  $t$ -th hop relations for the  $(t-1)$  hop relation path  $p_{t-1,i}$ .

After  $h$ -step prediction, we can obtain  $K^h$  relation paths. The score of the relation path  $p_{t,i}$  is the product of the probabilities of all relations in the path:

$$\text{Score}(p_{t,i}) = \text{Score}(r_{i,1} | r_{i,2} | \dots | r_{i,t}) = \prod_{l=1}^t d'_{r_{i,l}}, \quad i = 1, 2, \dots, K^h \quad (15)$$

We employ a training method similar to hop prediction. For the question  $q$  and the ground truth relation path, we concatenate question  $q$  and  $(t-1)$  hop relation path as the inputs  $Q_{t-1}$ . We encode  $Q_{t-1}$  into vector representation and feed it into the linear classification layer to obtain the  $t$ -th hop relation distribution  $D'_{r,t}$ . The ground truth distribution  $D_{r,t}$  is represented as a one-hot vector with a probability of 1 for the true  $t$ -th hop relation and a probability of 0 for other relations. We use the **cross-entropy loss** for parameter optimization.

**Triple Sampling.** We arrange the predicted relation paths in descending order of scores and sequentially sample reasoning paths composed of triples from KG until the number of reasoning paths reaches  $M$ . These reasoning paths are utilized as the relevant knowledge to augment LLMs on KGQA.

**4.1.2 Rewrite: KG-to-Text.** The core of our rewrite module is to transform structured triples into free-form text based on a KG-to-Text model. We start by training an open-source LLM based on question-related graph-text pairs. Given the graph  $G$  and corresponding free-form text  $y$ . We verbalize the triples in graph  $G$  into triple-form text  $x$  by concatenating the subject, relation, and object. Then, we transform triple-form text  $x$  to graph-to-text transformation prompt  $p_1$  via a template  $T_1$ : “Your task is to transform a knowledge graph to a sentence or multiple sentences. The knowledge graph is: {triple-form text  $x$ }. The sentence is:”. We take prompt  $p_1$  and free-form text  $y$  as the input and output respectively and adopt the teacher forcing strategy in training. Formally, given the prompt  $p_1$ , ground truth output sequence  $y = [y_1, y_2, \dots, y_T]$  and model vocabulary  $[v_1, v_2, \dots, v_V]$ , the model predicts the probability distribution  $D'_{v,t}$  of tokens at step  $t$  based on prompt  $p_1$  and previous  $(t-1)$  steps correct tokens  $y_1, y_2, \dots, y_{t-1}$ :

$$D'_{v,t} = [d'_{v_1}, d'_{v_2}, \dots, d'_{v_V}] \quad (16)$$

t 번째 토큰을 예측하는데 이전 토큰들의 정보가 쓰임

where  $d'_{v_c}$  is the probability of  $v_c$  given  $p_1$  and  $y_1, y_2, \dots, y_{t-1}$ :

$$d'_{v_c} = P(v_c | p_1, y_1, y_2, \dots, y_{t-1}), \quad c = 1, 2, \dots, V \quad (17)$$

Ground truth distribution  $D_{v,t}$  is a one-hot vector with a probability of 1 for the true next token  $y_t$  and probabilities of 0 for other tokens:

$$D_{v,t} = [d_{v_1}, d_{v_2}, \dots, d_{v_V}] \quad (18)$$

$$d_{v_c} = \begin{cases} 1, & v_c = y_t \\ 0, & v_c \neq y_t \end{cases} \quad c = 1, 2, \dots, V \quad (19)$$

The **cross entropy loss**  $L_{CE}$  is used to update the parameters:

$$J_t = -D_{v,t} \log D'_{v,t} = - \sum_{c=1}^V d_{v_c} \log d'_{v_c} \quad (20)$$

$$L_{CE} = \frac{1}{T} \sum_{t=1}^T J_t \quad (21)$$

While answering questions, each reasoning path is first linearized into triple-form text and then transformed into a prompt via template  $T_1$ . The prompt is fed into the fine-tuned LLMs to obtain the corresponding textual description. These descriptions are consolidated into a paragraph as the question-relevant knowledge to enhance the performance of the LLMs.

**4.1.3 Answer: Knowledge Text Enhanced Reasoning.** To integrate the generated knowledge  $y$  with the question  $q$ , we devise a template  $T_2$ : “Below are the facts that might be relevant to answer the question: {free-form text  $y$ } Question: {question  $q$ } Answer:”. We map the free-form text  $y$  and the question  $q$  to KG-augmented prompt  $p_2$  using template  $T_2$ . Then we input prompt  $p_2$  into the question-answering model and collect output as the predicted answer  $a$ .



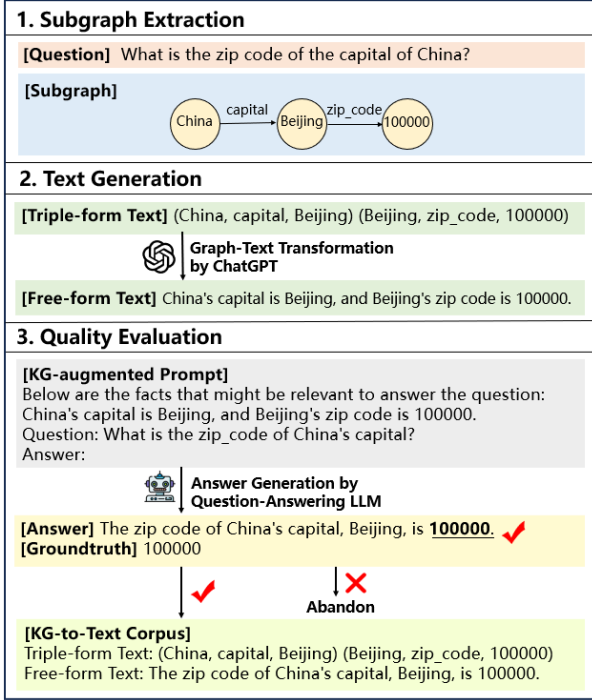


Figure 5: Our designed KG-to-Text corpus generation method has three steps: subgraph extraction, text generation, and quality evaluation.

## 4.2 Corpus Generation

Existing KGQA benchmarks do not provide graph-text pairs for the question-answering task. Therefore, we design a KGQA task-driven corpus generation method. Considering ChatGPT’s powerful natural language understanding and generation capabilities, inspired by [5, 16, 22], we adopt ChatGPT as the corpus generator. As shown in Figure 5, this process contains three steps: subgraph extraction, text generation, and quality evaluation.

**4.2.1 Subgraph Extraction.** For KGQA benchmarks that provide relation paths or reasoning triples (e.g. MetaQA [44], PathQuestion [45]), we obtain subgraphs by querying KG based on annotations. For benchmarks annotated with SPARQL (e.g. WebQuestionSP [42], ComplexWebQuestion [33]), we modify SPARQL query to retrieve intermediate entities and construct question-related subgraphs.

**4.2.2 Text Generation.** Given the question  $q$ , we first verbalize the related subgraph  $G$  into triple-form text  $x$  by concatenating the subject, relation, and object. Then, template  $T1$  is adopted to transform triple-form text  $x$  to graph-to-text transformation prompt  $p1$ . Finally, we input prompt  $p1$  into ChatGPT to obtain the corresponding free-form text  $y$ .

**4.2.3 Quality Evaluation.** Due to the absence of annotations, common evaluation metrics including BLEU, METEOR, and ROUGE cannot be employed. Considering the purpose of the generated text is to strengthen the performance of LLMs on KGQA, we assess the

quality of free-form text  $y$  based on the feedback from question-answering models. We map the free-form text  $y$  and question  $q$  to KG-augmented prompt  $p2$  via template  $T2$ . Then, we feed prompt  $p2$  to the question-answering model and get the predicted answer  $a$ . In light of the fact that LLMs typically provide textual passages as responses instead of single answer entities, we employ `hit@1` as a metric to assess the correctness of answer  $a$ . To put it simply, if answer  $a$  contains at least one answer entity, the question is considered answered correctly. In this scenario, we collect the triple-form text  $x$  and free-form text  $y$  as a graph-text pair.

## 5 EXPERIMENTS

### 5.1 Datasets

**MetaQA** [44] is a large-scale multi-hop KGQA benchmark in the movie domain. It provides a knowledge graph including 135k triples, 43k entities, and 9 relations. It contains more than 400k questions divided into MetaQA 1-hop, MetaQA 2-hop, and MetaQA 3-hop based on the hop number of the question. Each question is annotated with the head entity, answers, and entity categories involved in the reasoning path. In this experiment, we choose MetaQA 2-hop as our benchmark and the “vanilla” version of the questions, totaling 148,724 questions (118,980 train, 14,872 dev, 14,872 test). We collect the gold relation paths based on the provided entity categories since there is only one type of relation between two categories of entities.

**WebQuestionsSP (WebQSP)** [42] is a smaller scale KGQA benchmark with a larger scale KG. It provides SPARQL queries for WebQuestions [2] and filters out questions that are not answerable. The remaining 4,737 questions (3,098 train, 1,639 test) are either 1-hop or 2-hop questions with corresponding topic entities, inferential chains, and SPARQL queries. Following [30], we prune the KG to contain only relations mentioned in the questions and the triples within 2 hops of mentioned entities. The pruned KG includes 1.8 million entities, 627 relations, and 5.7 million triples.

**WebQuestions (WebQ)** [2] collects questions from web pages via Google Suggest API. Following [32], we use a subset of this benchmark. Our method requires annotations for relation paths or SPARQL queries, which are not provided. Therefore, we choose the questions that are provided with SPARQL queries in WebQSP for training and testing. This results in 4,737 questions (3,098 train, 1,639 test). We use the same KG as WebQSP.

**ZhejiangQA (ZJQA)** is a Chinese KGQA dataset of 20,491 questions provided by Zhejiang Lab. We divide these questions into trainset, devset, and testset (14,999 train, 2,147 dev, 3,345 test). The questions are either 1-hop or 2-hop questions primarily in the robotic domains. Each question is provided with a head entity, answers, and a gold relation path. It also provides a KG with more than 11k triples, 9k entities, and 39 relations.

### 5.2 Large Language Models

Our proposed framework has two modules based on LLMs: KG-to-Text and question-answering. We use Llama-2 (7B, 13B), Flan-T5 (3B) for KG-to-Text and Llama-2 (7B, 13B), T5 (0.8B, 3B, 11B), Flan-T5 (80M, 3B, 11B), T0 (3B, 11B), ChatGPT for question-answering.

**Llama-2** [37] is an updated version of Llama-1 [36], which is trained on a wide range of public online data sources. Among

**Table 1: Experimental results of our proposed framework and baselines on WebQSP/WebQ, where we report hit@1(%). We implement our KG-to-Text model based on Llama-2-chat (13B). The results of the baselines are copied from the original paper. The best score is emphasized in bold.**

Method	WebQSP					WebQ				
	T5 <sub>(0.8B)</sub>	T5 <sub>(3B)</sub>	T5 <sub>(11B)</sub>	T0 <sub>(3B)</sub>	T0 <sub>(11B)</sub>	Flan-T5 <sub>(80M)</sub>	Flan-T5 <sub>(3B)</sub>	Flan-T5 <sub>(11B)</sub>	T0 <sub>(3B)</sub>	T0 <sub>(11B)</sub>
KAPING	34.70	25.41	24.91	52.28	62.85	-	-	-	-	-
Sen et al.	-	-	-	-	-	45.52	55.58	59.79	53.33	55.64
<b>Ours</b>	<b>69.41</b>	<b>72.11</b>	<b>79.36</b>	<b>74.02</b>	<b>68.55</b>	<b>59.95</b>	<b>68.92</b>	<b>73.71</b>	<b>74.02</b>	<b>68.55</b>

different variants, we adopt Llama-2-chat (7B, 13B) as our KG-to-Text and question-answering model. For ZJQA, we employ the Chinese version of this model, Chinese-Alpaca-2 (7B, 13B)<sup>2</sup>.

T5 [24] is an encoder-decoder model pre-trained on multiple tasks in a text-to-text format. Following [1], we use the LM-adapted version<sup>3</sup> of T5 as the question-answering model on WebQSP to ensure a fair comparison.

Flan-T5 [4] is an extension of T5, which is further instruction tuned on a large-scale collection of automatically generated instructions from existing datasets. We use Flan-T5-XL (3B) as KG-to-Text model for MetaQA and Flan-T5-Small (80M), Flan-T5-XL (3B), Flan-T5-XXL (11B) as question-answering models for WebQ. We do not use this model for ZJQA since it does not support Chinese.

T0 [28] is fine-tuned from T5 over a variety of prompts to improve zero-shot generalization performance. We use T0 (3B, 11B) as the question-answering model to compare our proposed framework with previous work on WebQSP and WebQ.

ChatGPT<sup>4</sup> is a large language model built on GPT-3.5 developed by OpenAI. It is pre-trained on enormous corpus and human annotations and excels in understanding and generating human-like text. Specifically, we use GPT-3.5 Turbo in this experiment. We cannot fine-tune ChatGPT since it is not yet open source. Therefore, we access it via API<sup>5</sup> and use it as the question-answering model.

### 5.3 Baselines

We compare the proposed framework with previous KG-augmented LLM methods for KGQA on WebQSP and WebQ:

**KAPING** [1] verbalizes triples by concatenating subject, relation, and object. During the retrieval process, verbalized triples and questions are embedded into vector space with an off-the-shelf sentence embedding model. The vector similarity between questions and verbalized triples is calculated to retrieve triples as augmented knowledge. The retrieved triples and questions are concatenated in the form of the prompt and fed into the question-answering model to obtain answers. We compare with this baseline on WebQSP.

**Sen et al.** [32] improves the retrieval method of KAPING [1]. This work utilizes Rigel [27] to predict the distribution over relations for each hop and retrieve triples based on the calculated relation distribution and question entities. The retrieved triples are linear verbalized and concatenated with the question via prompt

and input into the question-answering model to generate the answer. We compare with this baseline on WebQ.

### 5.4 Evaluation Metrics

Following the evaluation settings in previous work on generative KGQA [29, 31, 43], we use hit@1, which measures whether the generated answer includes at least one answer entity. Both the predicted answers and the answer entities are converted to lowercase to mitigate matching issues arising from letter-case differences.

### 5.5 Implementation Details

**WebQSP/WebQ** We parse the SPARQL query to extract the gold relation path for each question. Bert-base-uncased is used as the classification model for hop prediction and relation path prediction. We modify the SPARQL query and obtain intermediate entities to construct the gold subgraph for each question. All the questions in train split are used for corpus generation. Over 12k graph-text pairs are generated and utilized for supervised fine-tuning. In addition, we exclude 11 test samples without answers.

**MetaQA** We skip the hop prediction step and use the provided ground truth hop number. Bert-base-uncased is used as the classification model for relation path prediction. We randomly sample 17k questions from train split for corpus generation and generate more than 13k graph-text pairs for KG-to-Text fine-tuning.

**ZJQA** We use bert-base-chinese for hop prediction and relation path prediction. We randomly sample 14k questions from train split for corpus generation and obtain over 13k KG-to-Text annotations.

For WebQSP and WebQ, we set K=5 (i.e. sample 5 relations as the next possible relations for each predicted relation path) and M=5 (i.e. sample up to 5 reasoning paths for each question) in the subgraph retrieval process. For MetaQA and ZJQA, we set K=3 and M=5. ChatGPT is used as the question-answering model in the corpus generation for all benchmarks. Our framework is implemented using Pytorch<sup>6</sup>, Transformers<sup>7</sup> and Peft libraries<sup>8</sup>. For efficient training, all fine-tuning processes adopt LoRA [8]. We train KG-to-Text models on 4 NVIDIA Tesla V100 GPUs for 10 epochs with a total batch size of 128 and run inference over KG-to-Text models and question-answering models on 1 NVIDIA Tesla V100 GPUs. We use AdamW optimizer with initialized learning rate 1e-4. The LoRA rank, LoRA alpha, and LoRA dropout are 64, 128, and 0.05.

<sup>2</sup><https://github.com/ymcui/Chinese-LLaMA-Alpaca-2>

<sup>3</sup>[https://github.com/google-research/text-to-text-transfer-transformer/blob/main/released\\_checkpoints.md](https://github.com/google-research/text-to-text-transfer-transformer/blob/main/released_checkpoints.md)

<sup>4</sup><https://openai.com/blog/chatgpt>

<sup>5</sup><https://api.openai.com/>

<sup>6</sup><https://pytorch.org/>

<sup>7</sup><https://huggingface.co/docs/transformers/main/index>

<sup>8</sup><https://huggingface.co/docs/peft/index>

**Table 2: Experimental results of different knowledge representation formats on MetaQA and ZJQA, where we report hit@1(%). The below section of the table is the results of our framework implemented by different KG-to-Text models. The number inside the parentheses behind the LLMs denotes its parameter size. The best score for each question-answering model is emphasized in bold.**

Knowledge Format	MetaQA			ZJQA		
	Llama-2-chat <sub>(7B)</sub>	Llama-2-chat <sub>(13B)</sub>	ChatGPT	Llama-2-chat <sub>(7B)</sub>	Llama-2-chat <sub>(13B)</sub>	ChatGPT
No Knowledge	34.31	34.08	60.58	14.95	15.75	15.25
Triple Knowledge	96.98	96.78	93.19	88.67	84.48	92.02
MTL Knowledge	86.32	87.18	80.66	-	-	-
MVP Knowledge	96.13	98.04	92.91	-	-	-
<b>Flan-T5<sub>(3B)</sub></b>	97.60	98.78	<b>97.71</b>	-	-	-
<b>Llama-2-chat<sub>(7B)</sub></b>	97.74	98.92	97.69	<b>93.09</b>	<b>93.21</b>	<b>92.91</b>
<b>Llama-2-chat<sub>(13B)</sub></b>	<b>97.77</b>	<b>99.07</b>	97.63	91.48	92.11	92.11

## 5.6 Main Results

Table 1 shows the overall results of our proposed framework and baselines on WebQSP and WebQ. Our KG-to-Text model is implemented based on Llama-2-chat (13B). For the question-answering model, we choose T5 (0.8B, 3B, 11B), T0 (3B, 11B) and Flan-T5 (80M, 3B, 11B), T0 (3B, 11B) as question-answering models on WebQSP and WebQ respectively. Experimental results demonstrate that our framework outperforms baselines by a large margin across various LLMs. Notably, it exhibits the most significant advantages on T5. This suggests that T5, which is pre-trained solely on textual data, may have limitations in comprehending structured data. This proves that **transforming triple-form text into free-form text can enable LLMs to better understand the provided factual knowledge and enhance their capabilities on KGQA.**

## 5.7 Ablation Study

We conduct a comparative analysis of different knowledge representation formats and implement our framework utilizing a range of LLMs on MetaQA and ZJQA datasets. This ablation study is performed to investigate the impact of knowledge representation formats and LLMs on KGQA.

Specifically, we use Llama-2-chat (7B, 13B), Flan-T5-XL (3B) as the KG-to-Text model and Llama-2-chat (7B, 13B), ChatGPT as the question-answering model. We compare knowledge generated by our framework with no knowledge, triple knowledge, and knowledge generated by the off-the-shelf KG-to-Text model:

**No Knowledge** Questions are directly fed into LLMs without additional knowledge. We set this knowledge representation format to explore how much improvement KG-augmented methods can bring to LLMs on KGQA.

**Triple Knowledge** Triple knowledge is the most common strategy used in previous work. We first deduplicate the retrieved triples to mitigate semantic redundancy and then simply verbalize each triple by concatenating the subject, relation, and object.

**MVP Knowledge** To verify the effectiveness of the fine-tuning process, we choose MVP [35] as the off-the-shelf KG-to-Text model. MVP is a text-generation LLM that is first pre-trained on 77 datasets over 11 diverse natural language generation (NLG) tasks in a supervised text-to-text format and then further pre-trained task-specific soft prompts to enhance the model’s ability on specific tasks. We

use MVP-data-to-text, a variant of MVP which is pre-trained on labeled data-to-text datasets to perform KG-to-Text in the zero-shot scenarios. We do not use this knowledge representation format for ZJQA since MVP does not support Chinese.

**MTL Knowledge** We select MTL-data-to-text as another off-the-shelf KG-to-Text model. This model is a different variant of MVP and is pre-trained on a mixture of labeled data-to-text datasets. However, compared with MVP-data-to-text, it lacks training on other NLG tasks and task-specific soft prompts pre-training. We do not use this model for ZJQA due to its lack of support for Chinese.

Table 2 shows the experimental results of different knowledge representation formats based on various LLMs as the question-answering models. Compared with other knowledge representation formats, knowledge generated by our framework shows further improvement across multiple LLMs on KGQA.

## 5.8 Analysis

**Impact of Knowledge Representation Formats** Table 2 demonstrates that different representations of the same retrieved triples have a significant influence on KGQA. The above section of Table 2 shows the results of the baseline representation formats. Without incorporating external knowledge, the performance of the question-answering model is the worst. This indicates that LLMs are incapable of storing all factual knowledge within their vast parameters, leading to issues of factual inaccuracy and knowledge missing. Other KG-augmented LLM methods outperform this baseline by a large margin, demonstrating the effectiveness of incorporating KG knowledge relevant to the questions. Among them, MTL knowledge exhibits the smallest improvement. This is because MTL is only pre-trained on multiple data-to-text datasets and its natural language understanding and generation capabilities are not strong enough. Therefore, the transformation from subgraphs to text loses semantic information and leads to limited enhancement. MVP-generated knowledge and triple-form knowledge result in comparable results in the question-answering task. The reason is that MVP is pre-trained on multiple NLG tasks and further fine-tuned on data-to-text datasets, possessing stronger capabilities in text comprehension and generation compared to MTL. Nevertheless, it lacks fine-tuning on domain-specific KG-to-Text corpus, leading to similar results as triplet knowledge. Triple-form knowledge is the most common approach in previous work. However, experimental

**Table 3: The detailed comparison results over no knowledge and triple knowledge. The below section of the table is the results of our framework implemented by different KG-to-Text models. We use “Helpful” and “Harmful” to evaluate the effectiveness of different knowledge representation formats and KG-to-Text methods. We adopt Llama-2-chat (13B) as the question-answering model. The best score is emphasized in bold.**

Knowledge Format	No Knowledge Baseline		Triple Knowledge Baseline	
	Helpful $\uparrow$	Harmful $\downarrow$	Helpful $\uparrow$	Harmful $\downarrow$
Triple Knowledge	9,428	103	-	-
MTL Knowledge	8,201	303	334	1,761
MVP Knowledge	9,583	70	443	255
<b>Flan-T5</b> <sub>(3B)</sub>	9,672	49	453	155
<b>Llama-2-chat</b> <sub>(7B)</sub>	9,690	46	456	137
<b>Llama-2-chat</b> <sub>(13B)</sub>	<b>9,703</b>	<b>38</b>	<b>459</b>	<b>119</b>

results show that LLMs still struggle to extract semantics effectively from triples. This suggests that **LLMs prefer receiving textual knowledge, as they have been pre-trained on massive corpora and structured triples are just a part of it.** The below section of Table 2 shows the results of our framework. Our framework employs multiple KG-to-Text models, surpassing all baselines on various question-answering models. This not only demonstrates that our KG-to-Text method can generate answer-sensitive textual knowledge but also highlights the robust applicability of our framework to mainstream LLMs.

**Impact of LLMs** The below section of Table 2 illustrates different performances of our proposed framework based on different LLMs. Generally, Llama-2-chat is better at KG-to-Text compared with Flan-T5-XL. We speculate this disparity can be attributed to differences in model parameters. The smaller size of Flan-T5-XL (3B) compared to Llama-2-chat (7B, 13B) may be responsible for the inferior performance. The two-parameter versions of Llama perform comparably. We believe this is due to the relative simplicity of KG-to-Text compared to other NLG tasks and 7B parameters are enough. For question-answering models, Llama-2-chat (13B) achieves the best performance on MetaQA while ChatGPT demonstrates the poorest performance. It is noteworthy that ChatGPT outperforms Llama-2-chat (7B, 13B) when no knowledge is provided. This observation indicates that **ChatGPT has a greater capacity to retain knowledge due to its significantly larger parameter size.** However, it does not leverage relevant knowledge as effectively as Llama-2-chat. For ZJQA, we employ the Chinese version of Llama-2-chat, Chinese-Alpaca-2. It is further trained based on Llama-2 using a substantial Chinese corpus and instruction data. Despite having significantly different model parameters, it achieves performance on ZJQA that matches or even surpasses ChatGPT. This emphasizes the significance of continued pre-training. Chinese-Alpaca-2 (13B) excels in almost all knowledge representation formats but lags behind ChatGPT in triple-form knowledge. Apart from differences in model parameters, we posit this is because ChatGPT has been trained on a wider range of structured data, enabling it to extract semantic information from structured data more effectively.

**Helpfulness/Harmfulness of Generated Knowledge** We further analyze the experimental results of MetaQA to investigate the positive and negative impact of knowledge generated by different methods on question-answering models. We establish two baselines: “no knowledge” and “triple knowledge”, and compare

other knowledge formats with these two baselines. We count the number of questions where the baseline answers incorrectly but other knowledge formats answer correctly (i.e., helpful), as well as the number of questions where the baseline answers correctly but other knowledge formats answer incorrectly (i.e., harmful). We choose Llama-2-chat (13B) as the question-answering model. The detailed information is presented in Table 3.

The comparison with no knowledge baseline indicates that the knowledge representation format employed in our framework can better assist LLMs on KGQA and has fewer adverse effects. Besides, experimental results on the triple knowledge baseline reveal that our KG-to-Text method can generate the most informative textual statements for KGQA. **This highlights the necessity of fine-tuning KG-to-Text models.**

## 6 CONCLUSION

In this paper, we propose Retrieve-Rewrite-Answer, a KG-to-Text enhanced LLMs framework for KGQA. Our framework adopts an answer-sensitive KG-to-Text method to generate textual knowledge which is most informative for KGQA. To address the challenge of missing annotation data of KG-to-Text, we design a KG-to-Text corpus generation method based on the feedback of question-answering models. Experimental results demonstrate that our framework outperforms previous KG-augmented methods by a large margin. Besides, compared to other verbalization methods, the free-form knowledge text generated by our framework is more comprehensible to LLMs and achieves superior results on KGQA. However, this framework requires fine-tuning KG-to-Text models on graph-to-text corpus and does not incorporate extra data sources besides KG. In future work, we plan to explore the integration of additional knowledge resources and design an approach capable of generating question-related knowledge in zero-shot scenarios.

## REFERENCES

- [1] Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-Augmented Language Model Prompting for Zero-Shot Knowledge Graph Question Answering. In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*. Association for Computational Linguistics, Toronto, Canada, 78–106. <https://doi.org/10.18653/v1/2023.nlrse-1.7>
- [2] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic Parsing on Freebase from Question-Answer Pairs. In *EMNLP. ACL*, 1533–1544.
- [3] Deng Cai and Wai Lam. 2020. Graph Transformer for Graph-to-Sequence Learning. In *AAAI*. AAAI Press, 7464–7471.
- [4] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al.



2022. Scaling Instruction-Finetuned Language Models. (2022). <https://doi.org/10.48550/arXiv.2210.11416> arXiv:2210.11416
- [5] Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen Xu, Wei Liu, Ninghao Liu, et al. 2023. AugGPT: Leveraging ChatGPT for Text Data Augmentation. (2023). <https://doi.org/10.48550/arXiv.2302.13007> arXiv:2302.13007
- [6] Zhijiang Guo, Yan Zhang, Zhiyang Teng, and Wei Lu. 2019. Densely Connected Graph Convolutional Networks for Graph-to-Sequence Learning. *Trans. Assoc. Comput. Linguistics* 7 (2019), 297–312.
- [7] Jiuzhou Han and Ehsan Shareghi. 2022. Self-supervised Graph Masking Pre-training for Graph-to-Text Generation. In *EMNLP*. Association for Computational Linguistics, 4845–4853.
- [8] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*. OpenReview.net.
- [9] Nan Hu, Yike Wu, Guilin Qi, Dehai Min, Jiaoyan Chen, Jeff Z Pan, and Zafar Ali. 2023. An empirical study of pre-trained language models in simple knowledge graph question answering. *World Wide Web* (2023), 1–32.
- [10] Quzhe Huang, Mingxu Tao, Zhenwei An, Chen Zhang, Cong Jiang, Zhibin Chen, Zirui Wu, and Yansong Feng. 2023. Lawyer LLaMA Technical Report. (2023). <https://doi.org/10.48550/arXiv.2305.15062> arXiv:2305.15062
- [11] Robert L. Logan IV, Nelson F. Liu, Matthew E. Peters, Matt Gardner, and Sameer Singh. 2019. Barack’s Wife Hillary: Using Knowledge Graphs for Fact-Aware Language Modeling. In *ACL (1)*. Association for Computational Linguistics, 5962–5971.
- [12] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of Hallucination in Natural Language Generation. *ACM Comput. Surv.* 55, 12 (2023), 248:1–248:38.
- [13] Pei Ke, Haozhe Ji, Yu Ran, Xin Cui, Liwei Wang, Linfeng Song, Xiaoyan Zhu, and Minlie Huang. 2021. JointGT: Graph-Text Joint Representation Learning for Text Generation from Knowledge Graphs. In *ACL/IJCNLP (Findings) (Findings of ACL, Vol. ACL/IJCNLP 2021)*. Association for Computational Linguistics, 2526–2538.
- [14] Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2022. Internet-Augmented Dialogue Generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 8460–8478. <https://doi.org/10.18653/v1/2022.acl-long.579>
- [15] Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text Generation from Knowledge Graphs with Graph Transformers. In *NAACL-HLT (1)*. Association for Computational Linguistics, 2284–2293.
- [16] Md Tahmid Rahman Laskar, Mizanur Rahman, Israt Jahan, Enamul Hoque, and Jimmy Huang. 2023. CQSumDP: A ChatGPT-Annotated Resource for Query-Focused Abstractive Summarization Based on Debatepedia. (2023). <https://doi.org/10.48550/arXiv.2305.06147> arXiv:2305.06147
- [17] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *ACL*. Association for Computational Linguistics, 7871–7880.
- [18] Jiacheng Liu, Skyler Hallinan, Ximing Lu, Pengfei He, Sean Welleck, Hannaneh Hajishirzi, and Yejin Choi. 2022. Rainier: Reinforced Knowledge Inspector for Commonsense Question Answering. In *EMNLP*. Association for Computational Linguistics, 8938–8958.
- [19] Dakuan Lu, Hengkui Wu, Jiaqing Liang, Yipei Xu, Qianyu He, Yipeng Geng, Mengkun Han, Yingsi Xin, and Yanghua Xiao. 2023. BBT-Fin: Comprehensive Construction of Chinese Financial Domain Pre-trained Language Model, Corpus and Benchmark. (2023). <https://doi.org/10.48550/arXiv.2302.09432> arXiv:2302.09432
- [20] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories. In *ACL (1)*. Association for Computational Linguistics, 9802–9822.
- [21] Diego Marcheggiani and Laura Perez-Beltrachini. 2018. Deep Graph Convolutional Encoders for Structured Data to Text Generation. In *INLG*. Association for Computational Linguistics, 1–9.
- [22] Lidia Ostyakova, Kseniia PetukhovaO, Veronika Smilga, and Dilyara ZharikovaO. 2023. Linguistic Annotation Generation with ChatGPT: a Synthetic Dataset of Speech Functions for Discourse Annotation of Casual Conversations. In *Proceedings of the International Conference “Dialogue*, Vol. 2023.
- [23] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. *Language Models are Unsupervised Multitask Learners*. Technical Report.
- [24] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* 21 (2020), 140:1–140:67.
- [25] Leonardo F. R. Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2021. Investigating Pretrained Language Models for Graph-to-Text Generation. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*. Association for Computational Linguistics, Online, 211–227. <https://doi.org/10.18653/v1/2021.nlp4convai-1.20>
- [26] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object Hallucination in Image Captioning. In *EMNLP*. Association for Computational Linguistics, 4035–4045.
- [27] Amir Saffari, Armin Oliya, Priyanka Sen, and Tom Ayoola. 2021. End-to-End Entity Resolution and Question Answering Using Differentiable Knowledge Graphs. In *EMNLP (1)*. Association for Computational Linguistics, 4193–4200.
- [28] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafei, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. 2022. Multitask Prompted Training Enables Zero-Shot Task Generalization. In *ICLR*. OpenReview.net.
- [29] Apoorv Saxena, Adrian Kochsiek, and Rainer Gemulla. 2022. Sequence-to-Sequence Knowledge Graph Completion and Question Answering. In *ACL (1)*. Association for Computational Linguistics, 2814–2828.
- [30] Apoorv Saxena, Aditya Tripathi, and Partha P. Talukdar. 2020. Improving Multi-hop Question Answering over Knowledge Graphs using Knowledge Base Embeddings. In *ACL*. Association for Computational Linguistics, 4498–4507.
- [31] Priyanka Sen, Alham Fikri Aji, and Amir Saffari. 2022. Mintaka: A Complex, Natural, and Multilingual Dataset for End-to-End Question Answering. In *COLING*. International Committee on Computational Linguistics, 1604–1619.
- [32] Priyanka Sen, Sandeep Mavadia, and Amir Saffari. 2023. Knowledge Graph-augmented Language Models for Complex Question Answering. In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*. Association for Computational Linguistics, Toronto, Canada, 1–8. <https://doi.org/10.18653/v1/2023.nlrse-1.1>
- [33] Alon Talmor and Jonathan Berant. 2018. The Web as a Knowledge-Base for Answering Complex Questions. In *NAACL-HLT*. Association for Computational Linguistics, 641–651.
- [34] Yiming Tan, Dehai Min, Yu Li, Wenbo Li, Nan Hu, Yongrui Chen, and Guilin Qi. 2023. Evaluation of ChatGPT as a Question Answering System for Answering Complex Questions. (2023). <https://doi.org/10.48550/arXiv.2303.07992> arXiv:2303.07992
- [35] Tianyi Tang, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. 2023. MVP: Multi-task Supervised Pre-training for Natural Language Generation. In *ACL (Findings)*. Association for Computational Linguistics, 8758–8794.
- [36] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. LLaMA: Open and Efficient Foundation Language Models. (2023). <https://doi.org/10.48550/arXiv.2302.13971> arXiv:2302.13971
- [37] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shriti Bhoale, et al. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. (2023). <https://doi.org/10.48550/arXiv.2307.09288> arXiv:2307.09288
- [38] Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaeckermann, Mohamed Amin, Pi-Chuan Chang, Andrew Carroll, Chuck Lau, Ryutaro Tanno, Ira Ktena, et al. 2023. Towards Generalist Biomedical AI. (2023). <https://doi.org/10.48550/arXiv.2307.14334> arXiv:2307.14334
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NIPS*. 5998–6008.
- [40] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *ICLR (Poster)*. OpenReview.net.
- [41] Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. 2023. HuaTuo: Tuning LLaMA Model with Chinese Medical Knowledge. (2023). <https://doi.org/10.48550/arXiv.2304.06975> arXiv:2304.06975
- [42] Wen-tau Yih, Matthew Richardson, Christopher Meek, Ming-Wei Chang, and Jina Suh. 2016. The Value of Semantic Parse Labeling for Knowledge Base Question Answering. In *ACL (2)*. The Association for Computer Linguistics.
- [43] Jun Yin, Xin Jiang, Zhengdong Lu, Lifeng Shang, Hang Li, and Xiaoming Li. 2016. Neural Generative Question Answering. In *IJCAI/AAAI Press*, 2972–2978.
- [44] Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander J. Smola, and Le Song. 2018. Variational Reasoning for Question Answering With Knowledge Graph. In *AAAI*. AAAI Press, 6069–6076.
- [45] Mantong Zhou, Minlie Huang, and Xiaoyan Zhu. 2018. An Interpretable Reasoning Network for Multi-Relation Question Answering. In *COLING*. Association for Computational Linguistics, 2010–2022.
- [46] Jie Zhu, Junhui Li, Muhua Zhu, Longhua Qian, Min Zhang, and Guodong Zhou. 2019. Modeling Graph Structure in Transformer for Better AMR-to-Text Generation. In *EMNLP/IJCNLP (1)*. Association for Computational Linguistics, 5458–5467.