

Large Language Models Can Better Understand Knowledge Graphs Than We Thought

Xinbang Dai^a, Yuncheng Hua^{b,*}, Tongtong Wu^c, Yang Sheng^d, Qiu Ji^d, Guilin Qi^a

^a*Southeast University, Nanjing, Jiangsu, China*

^b*The University of New South Wales, Sydney, New South Wales, Australia*

^c*Monash University, Melbourne, Victoria, Australia*

^d*Nanjing University of Posts and Telecommunications, Nanjing, Jiangsu, China*

Abstract

When we integrate factual knowledge from knowledge graphs (KGs) into large language models (LLMs) to enhance their performance, the cost of injection through training increases with the scale of the models. Consequently, there is significant interest in developing prompt strategies that effectively incorporate KG information into LLMs. However, the community has not yet comprehensively understood how LLMs process and interpret KG information in different input formats and organizations within prompts, and researchers often rely on trial and error. To address this gap, we design extensive experiments to empirically study LLMs' comprehension of different KG prompts. At the literal level, we reveal LLMs' preferences for various input formats (from linearized triples to fluent natural language text). At the attention distribution level, we discuss the underlying mechanisms driving these preferences. We then investigate how the organization of structured knowledge impacts LLMs and evaluate LLMs' robustness in processing and utilizing KG information in practical scenarios. Our experiments show that (1) linearized triples are more effective than fluent NL text in helping LLMs understand KG information and answer fact-intensive questions; (2) Different LLMs exhibit varying preferences for different organizational formats of triples; (3) LLMs with larger scales are more susceptible to noisy, incomplete subgraphs.

Keywords: Knowledge Graph, Large Language Model

*Corresponding author.

Email address: devin.hua@unsw.edu.au (Yuncheng Hua)

1. Introduction

Recent studies commonly utilize databases containing extensive factual knowledge, such as knowledge graphs (KGs), to reduce hallucination in language models and enhance the quality of their generated content [1]. In the era of pre-trained language models (PLMs), integrating KG knowledge within the model during the training process has garnered significant interest within the community [2, 3, 4, 5, 6]. However, as language models evolve, training large-scale language models (LLMs) with billions of parameters using KG data may encounter limitations, such as severe resource constraints or lack of public access to model architectures, training corpus, or training methods [7]. To address these issues, researchers increasingly focus on injecting external knowledge into LLMs through prompt engineering techniques [8, 9, 10, 11]. This lightweight approach feeds knowledge from external KGs into LLMs in the form of prompts, demonstrating its effectiveness in addressing various challenges that rely solely on factual knowledge. Some studies indicate that LLMs are highly sensitive to input patterns, and different input formats can impact model performance [12, 13, 14].

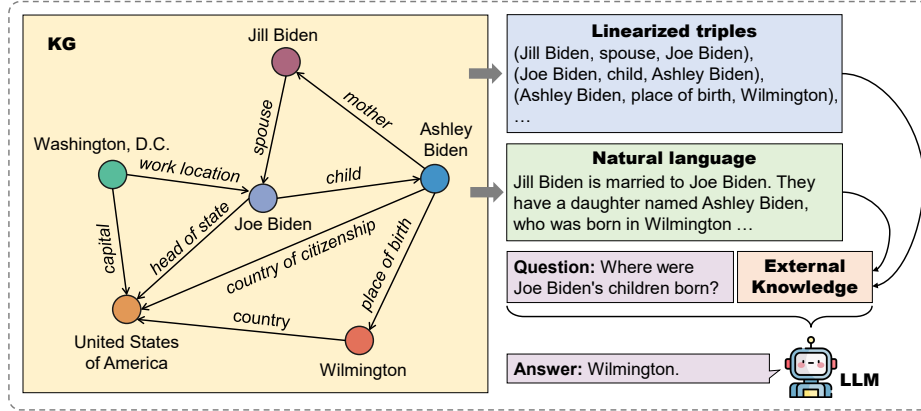


Figure 1: KG is processed into different input formats to provide LLM with knowledge.

As a highly structured form of knowledge, KGs can be integrated into LLMs in various ways. As illustrated in Figure 1, recent studies have explored directly feeding linearized triples from KGs into LLMs [15, 16]. Some other research employs KG-to-text generation approaches to convert structural knowledge prompts into natural language (NL) text, aiming to bridge the semantic gap between them [17, 18, 19, 20]. However, generating text from KGs becomes a significant challenge when dealing with multi-relation subgraphs containing numerous

triples (tens or even hundreds) [21]. Additionally, it is essential to recognize the necessity of bridging the semantic gap by rewriting linearized triples into NL text, as this process is complex and may be resource-intensive. The two types of approaches in Figure 1 highlight the community’s limited understanding of how LLMs comprehend KGs. To further assist researchers in constructing more efficient KG-related prompt strategies, it is crucial to understand which KG input format is most beneficial for prompting LLMs in addressing fact-intensive queries.

In this study, we aim to explore which KG input formats can better prompt LLMs [22]. We curate our question-answering (QA) evaluation data based on multiple datasets for KG-related tasks (such as KGQA and relation extraction) and ask LLMs to answer complex questions based on different KG input formats. Compared to other KG-related tasks (such as link prediction and KG-to-text generation), the QA task is more complicated and offers a unique evaluation advantage. Our questions encompass entity enumeration, counting, ranking, comparison, and truthfulness assessment, thereby thoroughly evaluating LLMs’ ability to adopt externally injected KG. LLMs may need to understand, retrieve, or associate external knowledge with their internal knowledge to better answer questions. Despite the complexity of these questions, the QA evaluation method is simple but relatively objective. We can accurately evaluate the performance of LLMs by directly comparing the predicted answers with the golden answers.

Specifically, we first evaluate the preferences of LLMs for different KG input formats at both the literal and attention levels. The KG input formats range from unordered triples to fluent NL. At the literal level, we design two complementary experiments, *Triple-to-Text* and *Text-to-Triple*, to assess LLMs’ preferences for KG input formats of varying scales and complexities. Compared to fluent NL text, we find that unordered linearized triples are more effective in assisting LLMs in answering knowledge-intensive questions. At the attention level, we analyze the attention distribution of LLMs across different input formats to investigate the reasons behind this phenomenon. Furthermore, we investigate the impact of the organization of linearized triples on LLM performance by employing various prompt strategies. Finally, we evaluate the robustness of LLMs to noisy and incomplete external subgraphs. We hope our findings can provide the community with insights for better designing KG-related prompting strategies to enhance LLM performance. Our main findings are as follows:

- *When using external knowledge to answer fact-intensive questions, LLMs prefer unordered structural data over fluent NL text.*
- *Sorting unordered linearized triples may not be necessary, as different LLMs*

prefer different prompt strategies. Researchers need to conduct meticulous experiments to develop more universal KG knowledge injection prompts.

- *In the input format of linearized triples, LLMs with larger parameters are less robust to such inputs, i.e., they are more susceptible to noisy or incomplete KGs, resulting in performance degradation.*

2. Related Works

2.1. Injecting KG into LLM During Training.

Injecting knowledge from KGs into LLMs during training has been extensively researched. This approach enables LLMs to grasp the semantics of KG embeddings through collaborative training [2, 3, 23, 4, 24, 6]. Although these methods have shown progress in smaller PLMs, their applicability to larger-scale LLMs presents challenges that require careful consideration of model architecture, training methods, and other aspects. In addition, injecting knowledge within training may result in insufficient and incorrect internalized knowledge in LLMs [25].

2.2. Integrating KG for LLM Using Prompt

Integrating external structural knowledge into prompts to enhance LLM capabilities has become a common strategy. Some studies [15, 16] directly provide LLMs with linearized structural knowledge as part of the prompts. [26] linearizes structural data into unified table rows, feeding them into the LLM to generate answers based on contextual examples. StructGPT [27] first selects entity and relation candidates, then uses multiple rounds of retrieving interface calls and LLMs as rankers to obtain answers.

However, other research suggests that inputs in NL format may be more model-friendly. These studies first convert structural knowledge into NL and then use it as prompts. UDT-QA [17] treats structural data as a form of knowledge expansion, converting it into NL text and adding it to the document repository for retrieval. [18] transforms subgraphs extracted from SPARQL queries into NL paragraphs, incorporating them into prompts to drive LLMs in generating NL questions. [19] believes that converting structural knowledge into high-quality NL text can substantially reduce the semantic gap between them. They employ KG-to-text generation models to rewrite structural knowledge and use it to help LLMs answer questions. KnowledgeNavigator [20] performs efficient reasoning on KGs and uses templates to convert structural reasoning paths into NL, guiding LLM reasoning. These methods assume that NL formats are more suitable for LLMs, while neglecting the potential information noise and loss incurred during the format conversion process.

3. Unveiling LLMs’ Understanding of KG in Different Input Format

In this section, we conduct a comprehensive study to determine which input format of KGs is more beneficial for prompting LLMs. Our analysis is conducted on two levels: literal and attentional. At the literal level, we observe phenomena, while at the attention distribution level, we explore the underlying reasons.

At the literal level (Section 3.1), we design two pipelines for prompting LLMs: *Triple-to-Text* and *Text-to-Triple*, and use a QA task to evaluate these two pipelines. In *Triple-to-Text*, we gradually convert the KG subgraph from the linearized triples into NL text generated by the model, which tests the LLM’s understanding of KGs across different scales and input formats when answering questions. The experiments indicate that prompts composed of unordered triples are more effective for LLMs when answering fact-intensive questions. However, this raises a concern: is the decline in LLM performance on NL texts due to potential quality issues in the text generated by the model? To address concerns about *Triple-to-Text*, we design a complementary pipeline: *Text-to-Triple*. Here, the NL text is manually written, and the corresponding subgraph is manually annotated, ensuring the quality of the KG provided to the LLM as a prompt. Since both the NL text and the KG are already given, we do not need to rely on generative models to convert the KG into NL text. This allows for a more fair and controlled experimental comparison between the NL text and the KG.

At the attention level (Section 3.2), we capture how the LLMs place greater focus on KG knowledge within the prompt. The experiments reveal that, regardless of whether a single KG format (linearized triples **or** NL text) or double formats (linearized triples **and** NL text) are provided, LLMs consistently demonstrate a greater capacity to capture answer information from linearized triples.

3.1. Literal Level Evaluation of LLM’s Understanding of KG

3.1.1. LLM’s Understanding of Triple-to-Text

In terms of the *Triple-to-Text* pipeline, we analyze the LLM’s ability to understand knowledge graph subgraphs from two perspectives: **KG subgraph scale** and **KG input format**. To quantitatively study the impact of the subgraph scale on LLMs, we propose *Core Reasoning Path Generation* and *Controlled Neighbour Node Search* to control the subgraph scale. Subsequently, the subgraphs are fed to the LLMs in various input formats to examine their preferences.

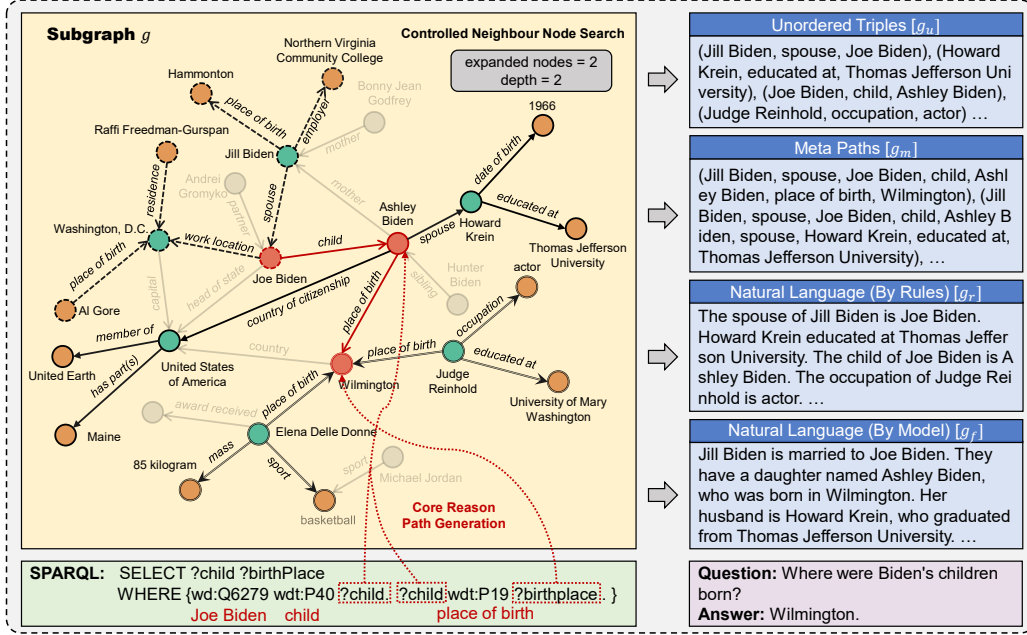


Figure 2: There are six categories of our expansion method. (1) *expanded nodes = 0, depth = 0*: only providing core reasoning paths; (2) *expanded nodes = 0.5, depth = 1*: expanding each node on each core path by one neighbouring node, with a 50% probability of deleting this expansion node; (3) *expanded nodes = 1, depth = 1*: expanding each node on each core path by one neighbouring node; (4) *expanded nodes = 2, depth = 1*: expanding each node on each core path by two neighbouring nodes; (5) *expanded nodes = 1, depth = 2*: starting from nodes on the core path, expanding to 2-hop neighbouring nodes, expanding one node at a time; (6) *expanded nodes = 2, depth = 2*: starting from nodes on the core path, expanding to 2-hop neighbouring nodes, expanding two nodes at a time (shown in this figure).

Core Reasoning Path Generation. Firstly, to obtain the core reasoning path C for each question, we utilize the SPARQL queries¹ provided in the KGQA dataset to retrieve all answers and their corresponding multiple reasoning paths. Specifically, we use the Wikidata endpoint to extract all constraint variable IDs from the SPARQL query for each question and then convert these IDs into their corresponding English labels. As shown in Figure 2, by arranging these constraint variables in sequence, we derive the core reasoning path in the KG (in red).

¹https://www.wikidata.org/wiki/Wikidata:SPARQL_tutorial

Controlled Neighbor Node Search. Once the reasoning path C is obtained, we control the number of adjacent nodes and the number of rounds for subgraph g expansion using the parameters *expanded nodes* and *depth*. As illustrated in Figure 2, when *expanded nodes* = 2 and *depth* = 2, two neighbouring nodes are randomly selected for expansion from each node on the reasoning path, and the expansion round is two. For example, *Joe Biden* expands to *Jill Biden* and *Washington, D.C.*, while *Ashley Biden* expands to *United States of America* and *Howard Krein*. The nodes from the first expansion round are shown in green, and those from the second round are in orange. As shown in Figure 2, the parameter *expanded nodes* is selected as {0.5, 1, 2}, and *depth* is selected as {1, 2}. Consequently, we construct six different subgraphs to evaluate the reasoning performance of the LLM.

KG Input Format. To investigate the impact of different knowledge formats for LLMs, we refer to some intermediate steps in recent works and devise five levels of KG injection methods: (1) Omitting Subgraph: we ask LLMs to respond to the questions without subgraph g , denoted as g_o . (2) Unordered Triples: we randomly shuffle all triples in subgraph g , denoted as g_u . (3) Meta Path: we connect triples sharing the same head or tail entity to construct meta-paths [28], denoted as g_m . (4) Natural Language (rule-based): we use heuristic rules to convert triples into NL text [29], denoted as g_r . (5) Natural Language (model-based): we employ a data-to-text generation model MVP [30] to convert g into fluent NL text [30], denoted as g_f . Figure 2 illustrates all input formats. There are a total of 25 possible combinations of extension and injection modes. Omitting subgraphs g_o is considered one mode because it does not incorporate subgraph knowledge.

Datasets. Wikidata [31] is a large-scale, high-quality KG that is frequently updated. We selected three KGQA datasets based on Wikidata, i.e., QALD-7 [32], LC-QuAD 2.0 [33], and KQAPro [34]. All of these three datasets already provide ground-truth SPARQL query annotations. These datasets serve as the foundation for generating our datasets. QALD-7 contains 215 training questions and 50 test questions. LC-QuAD 2.0 comprises 24k training questions and 6046 test questions. KQAPro includes 94k training questions and 10k test questions. From these datasets, we use a SPARQL endpoint to retrieve answers from Wikidata, filtering out questions with incorrect or unanswerable results. We also delete the questions where the core reasoning path cannot be extended to two rounds, such as when the path contains numerical or other attribute information. After filtering, QALD-7 retains 64 questions. For LC-QuAD 2.0 and KQAPro, considering the cost of calling the API, we randomly select 2000 questions from each dataset.

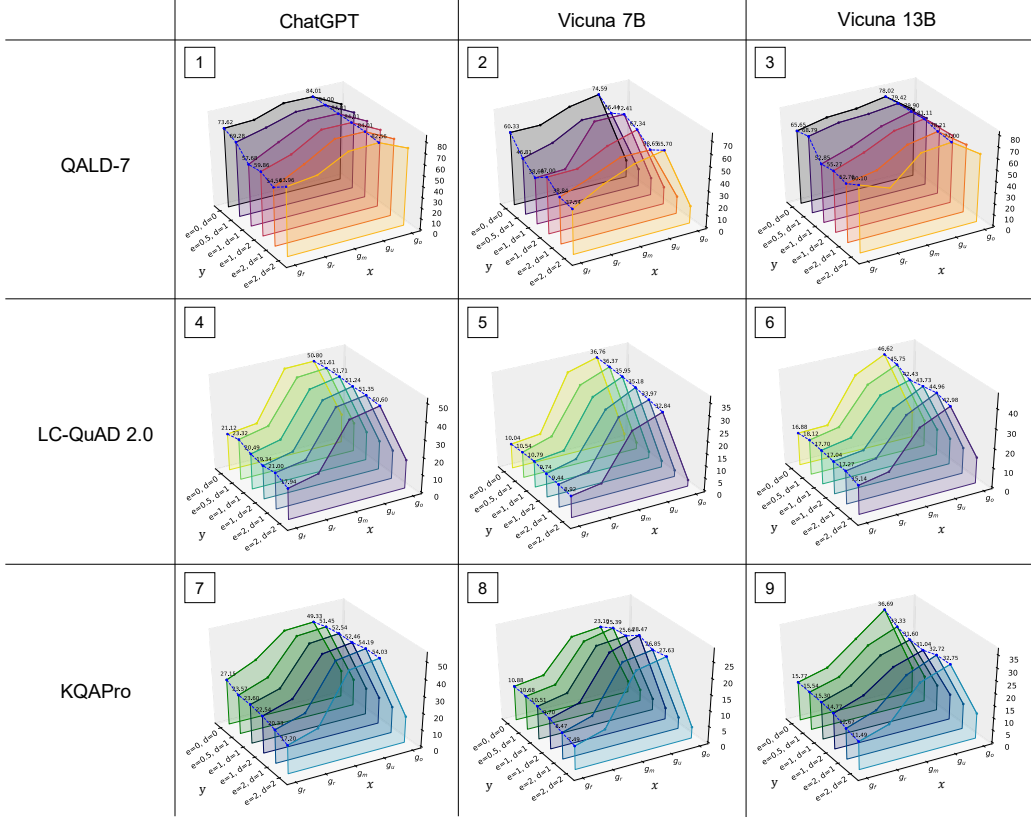


Figure 3: Performance of LLM in Triple-to-Text Pipeline. The x-axis represents various input formats, while the y-axis indicates different subgraph sizes. The parameter e denotes *expanded nodes*, and d represents *depth*.

Experimental Setup. We employ ChatGPT², Vicuna 7B and 13B [35] to evaluate the data, with all model parameters fixed. ChatGPT and Vicuna represent two mainstream series of LLMs, both demonstrating decent performance. To ensure a fair comparison, we use the exact match (EM) metric based on the standard SQuAD [36] to evaluate the performance of the LLM. Following the approach in [37], we maintain an alias table for entities in our answers, derived from Wikidata, to match different mentions of the same entity in the answers.

²<https://openai.com/index/chatgpt>

Results analysis. Figure 3 shows that the unordered linearized triple input g_u consistently outperforms other methods. The limited size (only 64 questions) of the QALD-7 dataset initially obscures the advantages of our process. However, as the questions’ complexity and the data’s scale increase (as demonstrated in LC-QuAD 2.0 and KQAPro), the linearized triples demonstrate better performance in the knowledge prompt. Moreover, the downward trend from g_u to g_f on the x-axis is consistent across different models and datasets, indicating that unordered triple knowledge can help LLM better answer multi-hop fact-intensive questions.

When examining the impact of subgraph size on LLM performance along the y-axis, we observe that LLMs do not experience significant performance degradation when faced with larger subgraphs. The blue dashed lines indicate the performance variation of LLMs under different subgraph scales with the input format of g_u and g_f . This decline is not as pronounced as the performance decline caused by changing input formats. This indicates that changes in the scale and structure of subgraphs have a limited impact on the performance of LLMs, while the input format is a more significant factor in prompts affecting LLMs. To further explore the preference of LLMs for linearized triple formats, in Section 4.2, we conduct experiments at the attention layer level. We observe that LLMs are more adept at focusing on knowledge relevant to the answers within linearized triples.

3.1.2. LLM’s Understanding of Text-to-Triple

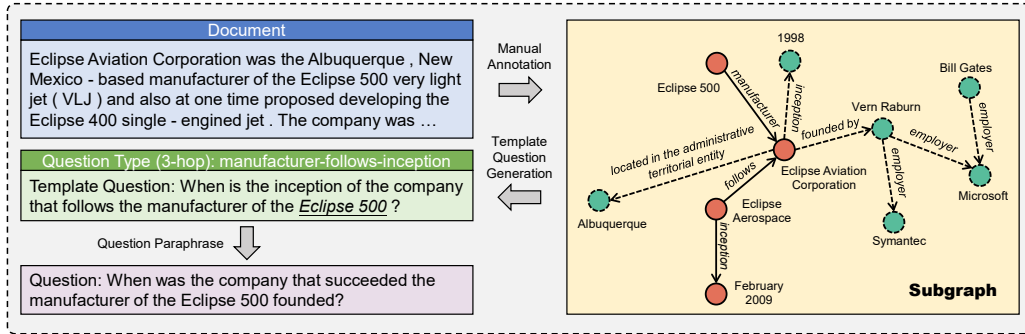


Figure 4: We employ a human-annotated KG mapping with a document for generating multi-hop questions, then evaluate the performance of LLM in answering these fact-related questions with *Documents* and *Subgraph*.

In *Text-to-Triple* pipeline, we use manually annotated NL text, rather than model-generated text, to evaluate LLMs’ ability to understand external knowledge. In *Triple-to-Text*, the quality of NL text g_f generated from KG may contain errors,

potentially affecting the ability of LLMs to understand KGs in NL format. In order to eliminate such concern, it is essential to provide LLMs with human-written NL text and establish a mapping from text to KG. The document-level relation extraction datasets, DocRED [38], derived from Wikidata, and ChemDisGene [39], based on the Comparative Toxicogenomics Database³ for biomedical data, provide human-crafted annotations for triples, which serve as ground-truth mappings for NL text. We process these datasets and generate two new QA datasets, addressing the challenge of ensuring NL text quality in triple generation. We also utilize the QA dataset derived from ChemDisGene to examine LLMs’ understanding capabilities when faced with different input formats of **domain-specific KG prompts**.

All documents in DocRED are human-written, and all mapped triples are manually annotated and aligned with entities and relations in Wikidata. The ChemDisGene document set is comprised of biomedical paper abstracts, with mapped triples manually curated by a team of biologists. In both datasets, the triples within each document form a complete, small-scale KG subgraph. We consider this subgraph as a complete structural representation of all entities and relations involved in the document. As shown in Figure 4, based on these subgraphs, we can extract all reasoning paths of different hops to generate fact-related questions.

Data Generation. To generate high-quality QA pairs from each document, we refer to the method used in MetaQA [40] for constructing multi-hop questions from triples. As shown in Figure 4, we fill the entities into various manually crafted multi-hop templates to ensure semantic coherence. Then, we use ChatGPT to paraphrase these questions, enhancing their diversity. DocRED consists of 5,053 Wikipedia documents, each associated with an annotated subgraph. We randomly select 800 documents that contain at least 3-hop paths and create 1-hop, 2-hop, and 3-hop questions to evaluate the LLM’s ability based on subgraphs. ChemDisGene consists of 523 abstracts, each annotated with the corresponding triples by the biologists. Based on the scale of these subgraphs corresponding to the documents, we generate 518 1-hop questions, 491 2-hop questions, and 243 3-hop questions. We separately provide the unordered triples and human-written NL documents to the LLM to answer these questions.

Experimental Setup. To eliminate doubts about the quality of NL knowledge generated by the model, we use the same LLM configuration and evaluation

³<https://ctdbase.org>

metric as *Triple-to-Text*, with the only difference being the replacement of g_f type knowledge with human-written documents.

LLMs	DocRED						ChemDisGene					
	1 hop		2 hop		3 hop		1 hop		2 hop		3 hop	
	Text	Triple	Text	Triple	Text	Triple	Text	Triple	Text	Triple	Text	Triple
ChatGPT	55.25	73.38	14.25	19.88	14.00	18.25	43.47	57.54	7.05	9.80	5.70	7.57
Vicuna 7B	34.88	50.13	9.50	11.00	8.63	10.50	25.21	32.53	3.49	7.28	2.72	4.91
Vicuna 13B	47.62	73.26	15.37	16.38	13.87	14.75	38.05	45.55	5.10	8.33	3.74	5.04
GPT-4	59.16	75.48	19.37	22.96	17.68	25.27	46.69	60.89	10.52	14.66	8.06	19.71
Llama3 13B	50.21	73.97	18.61	21.76	15.42	17.61	42.12	47.26	5.84	9.13	4.01	6.16

Table 1: Performance of LLMs in *Text-to-Triple* Pipeline. *Text* input is the document, and *Triple* input is randomly shuffled linearized triples.

Results Analysis. The data in Table 1 indicate that, compared to NL text, LLMs achieve significant improvements when using linearized triples, a finding consistent with those from *Triple-to-Text*. Additionally, we observe that for multi-hop, fact-intensive questions, knowledge prompts in text format may hinder LLMs from providing better answers. When addressing multi-hop questions, text format inputs require LLMs to perform reasoning across sentences or paragraphs, where unrelated information can adversely affect LLM performance. For instance, fluent NL texts may introduce noise (such as function words and determiners), hindering LLMs’ ability to recognize the core knowledge it should consider. LLMs do not struggle with understanding unordered linearized triples; on the contrary, this straightforward prompt type is more conducive to answering multi-hop questions.

In the domain-specific KG, compared to the general KG (Wikidata), the performance of LLMs shows a noticeable decline. This decline can be attributed to the fact that specialized domain knowledge is often not included in the pretraining corpus of LLMs. Consequently, LLMs lack the relevant expertise in specialized fields, which may result in a weaker grasp of domain-specific knowledge compared to general world knowledge. Notably, even when external knowledge is provided, enabling LLMs to understand and effectively utilize domain-specific information to answer questions remains a significant challenge.

Based on the complementary experiments of *Triple-to-Text* and *Text-to-Triple*, we draw an important conclusion: LLMs tend to prefer unordered linearized triples as the input format of KG.

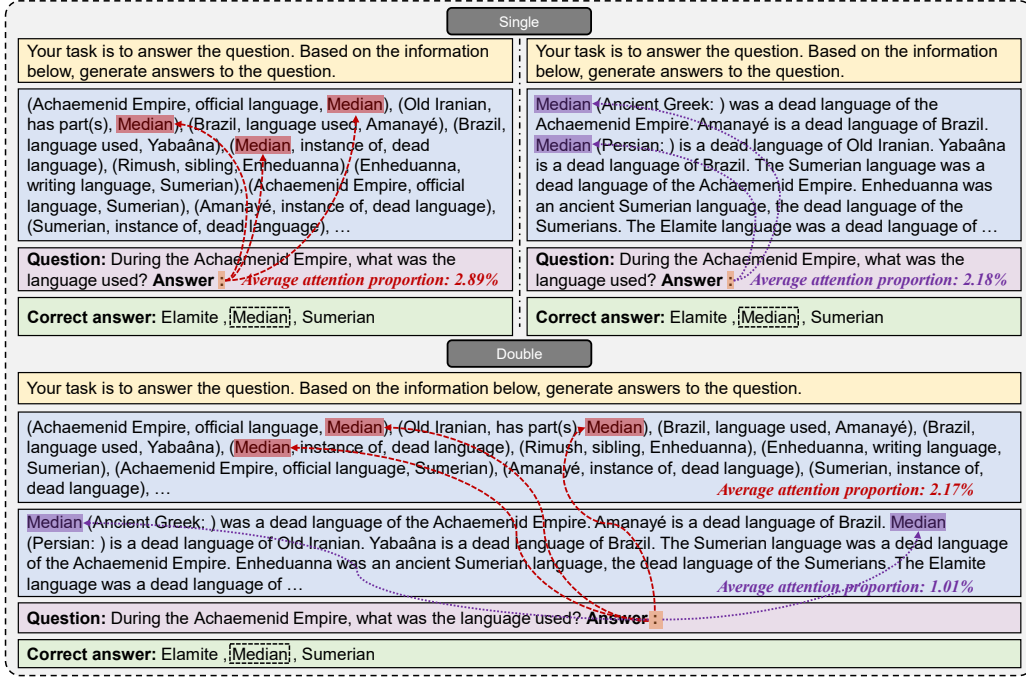


Figure 5: We examine the average attention proportion between the predicted labels (indicated by a colon ":") and the answer words (e.g., "Median"). Whether in the Single mode (providing knowledge in one format) or the Double mode (providing knowledge in both formats simultaneously), the LLM consistently pays more attention to the answers in the linearized triple format.

3.2. Attention Level Evaluation of LLM's Understanding of KG

This section explores how LLMs utilize KG information in different input formats, further revealing their ability to understand KGs. Specifically, we observe the attention distribution of LLMs towards answers in different KG input formats. This distribution illustrates the information's prioritized degree, thereby assisting us in analyzing which input format of KG is more beneficial for prompting LLMs. We introduce our evaluation method and then employ a quantitative analysis based on two datasets.

3.2.1. Attention Level Evaluation Method

As shown in Figure 5, the input format of knowledge is divided into unordered linearized triples g_u and NL text g_f . We employ two fair comparison modes for knowledge prompts: (1) Single: providing the LLM with only one format of knowledge (either g_u or g_f); (2) Double: providing the LLM with both formats of knowledge simultaneously (g_u and g_f).

During inference, the last token of the input, p , is used to trigger the model’s predicted label [41]. We examine the proportion of p ’s attention directed toward the tokens within the input that belong to the answer. This reflects the degree to which the LLM focuses on the answer within the provided knowledge. A higher proportion suggests a greater likelihood that the model will include the answer in its generated response. Considering that the number of answer occurrences may vary across different knowledge formats, we calculate the average attention proportion for each answer to ensure fairness. Therefore, we calculate the average attention proportion between the last token (the colon ":") and the answer word (for instance, "Median") as follows:

$$\overline{Att} = \frac{1}{n} \sum_{i=1}^n Att_{ans_i,p} \quad (1)$$

Here, $Att_{ans_i,p}$ represents the attention proportion between all tokens of answer i in the prompt and the prediction label p , and n denotes the total occurrences of answer tokens in the prompt.

3.2.2. Experimental Setup

Datasets and Metric. Our study selects two datasets from the *Triple-to-Text* experiment: LC-QuAD 2.0 and KQAPro. We utilize subgraphs with parameters *expanded_nodes*=2 and *depth*=2. The input modes are unordered linearized triples g_u and fluent NL text g_f . For a dataset D containing m questions, the attention proportion for question q_k ($k = 1, \dots, m$) is denoted as \overline{Att}_k . Consequently, the attention score of the LLM for the answers within the dataset D is calculated as:

$$Att_D = \frac{1}{m} \sum_{k=1}^m \overline{Att}_k \quad (2)$$

LLMs. We use the open-source Vicuna 7B and 13B models to examine how LLMs of different scales understand subgraphs with two input formats. In the final layer of the model, we sum the attention between each token and the prediction token p across all heads [41]. Then, we normalize all values to obtain each token’s attention proportion of p .

Results Analysis. The experimental results, as shown in Table 2, indicate that whether the two types of knowledge are provided simultaneously or separately, the model consistently shows a higher attention ratio to linearized triples. This suggests that LLMs can adapt to knowledge prompts that are less similar to the

	Vicuna 7B				Vicuna 13B			
	KQAPro		LC-QuAD 2.0		KQAPro		LC-QuAD 2.0	
	g_f	g_u	g_f	g_u	g_f	g_u	g_f	g_u
Single	2.46	3.80	2.08	2.67	2.78	4.69	2.10	3.78
Double	1.12	1.31	0.89	2.38	1.62	2.72	1.40	2.57

Table 2: The average attention proportion of LLM to the answer tokens in the prompt.

NL input format. Furthermore, when addressing fact-intensive questions, the LLM can identify the critical information and retrieve answers from linearized triple prompts. This explains why the LLM demonstrates greater interest in unordered linearized triples than NL text.

4. Comparing LLM’s Performance on Different KG Prompt Strategies

Based on the findings in Section 3, we observe that unordered linearized triples used as prompts yield more benefits for LLMs than fluent NL text. However, these triples are randomly ordered in prompts and lack a specific organization. Therefore, discussing the organization of linearized triples is essential for designing LLM-friendly prompts. In this section, we design various prompt strategies for linearized triples based on the *question relevance score*. We discover that different versions of LLMs exhibit distinct preferences for these strategies, highlighting the need to explore universal prompting techniques.

4.1. Question Relevance Score

4.1.1. Cross-encoder Based Scorer

To establish the criteria for sorting triples, we use the relevance score between the triples and the question. We construct a cross-encoder based on the BERT-base [42] model, similar to those described in [43] and [44]. The input $\tau_{q,t}$ consists of the concatenation of the question and the triple, enabling the model to develop deep cross-attention between them. Formally, we denote the joint embedding of the question and triple as $v_{q,t}$:

$$v_{q,t} = \text{red}(\text{BERT}_{\text{cross}}(\tau_{q,t})) \quad (3)$$

Here, $\tau_{q,t}$ represents the input representation of the question and triple, $\text{BERT}_{\text{cross}}$ is the BERT-based cross-encoder, and the function $\text{red}(\cdot)$ reduces the sequence of vectors produced by the encoder into one vector. Based on the experiments in [44],

we select the $\text{red}(\cdot)$ function as the last layer of the output of the [CLS] token. To score candidate entities, we apply a linear layer W to the embedding $v_{q,t}$:

$$s_{\text{cross}}(q, t) = v_{q,t}W \quad (4)$$

4.1.2. Experimental Setup

The scorer is based on the BERT-base model. We divide all the triples and questions from the *Text-to-Triple* dataset into the training set and a test set with an 8:2 ratio. Triples in the reasoning path linked to the question are labelled positive examples; otherwise, they are designated negative examples. For the cross-encoder, the batch size is 50; we experiment with initial learning rates of $\{5\text{e-}4, 2\text{e-}5, 5\text{e-}5, 2\text{e-}5\}$, and the learning rate decays every 3 epochs. We set the multiplicative factor, gamma, to update the learning rate to 0.2. Upon training, the model exhibits an accuracy of 98.89% in determining whether triples are related to the question, i.e., whether they are part of the core reasoning path.

4.2. Prompt Strategies

We perform the following operations on these triples using $s_{\text{cross}}(q, t)$: (1) Grouping: By setting thresholds at 0.3 and 0.8, we categorize the triples into high, medium, and low relevance groups relative to the question. (2) Ranking: We sort the candidate triples in descending order according to $s_{\text{cross}}(q, t)$. (3) Scoring: We append $s_{\text{cross}}(q, t)$ to each triple to indicate the relevance between the question and the triples to the LLM.

Finally, we input the triples and questions into the LLM to generate the answers. Table 3 illustrates three prompting strategies for the question: *What business structure did Frank Gehry design?* For grouping, we require the LLM to focus on higher relevant triples, thereby aiding in narrowing the search scope. For ranking, we expect the LLM to prioritize the foremost information in a sequence of triples. For scoring, we hope that the score will assist LLMs in retrieving relevant triples.

4.3. Datasets and Metric

We utilize all datasets from the *Triple-to-Text* experiment. We still employ subgraphs with parameters set to $\text{expanded_nodes}=2$ and $\text{depth}=2$. The originally unordered linearized triples are combined using three prompt strategies. The evaluation metrics follow the experimental setup in Section 3.1.1.

Strategies	Prompt
Grouping	<p>You are a QA assistant. For question: <i>What business structure did Frank Gehry design?</i> Refer to the following knowledge to response.</p> <p>Here are some triples that are highly relevant to the question: (DZ Bank building, architect, Frank Gehry), (Gehry Tower, instance of, office building), ...</p> <p>Here are some triples that are likely relevant to the question: (IAC Building, architect, Frank Gehry), (Gehry Tower, architect, Frank Gehry) ...</p> <p>Here are some triples that are less relevant to the question: (Toledo Museum of Art, architect, Frank Gehry), (Vlado Miluni, notable work, Dancing House), ...</p> <p>Answer:</p>
Ranking	<p>You are a QA assistant. For question: <i>What business structure did Frank Gehry design?</i> Refer to the following knowledge to response.</p> <p>These triples are sorted from high to low according to their relevance score to the question: (DZ Bank building, architect, Frank Gehry), (Dancing House, instance of, office building), (Gehry Tower, architect, Frank Gehry), (Dancing House, architect, Frank Gehry), (IAC Building, instance of, office building), ...</p> <p>Answer:</p>
Scoring	<p>You are a QA assistant. For question: <i>What business structure did Frank Gehry design?</i> Refer to the following knowledge to response.</p> <p>Each triple is followed by a relevance score to the question, which helps in solving the question: (DZ Bank building, architect, Frank Gehry) 0.9981, (Toledo Museum of Art, architect, Frank Gehry) 0.0019, Gehry Tower, instance of, office building) 0.998, (Vlado Miluni, notable work, Dancing House) 0.0023...</p> <p>Answer:</p>

Table 3: Different prompt strategies for triples based on relevance scores.

4.4. Results Analysis

The experimental results in Table 4 indicate that, compared to the unordered linearized triplet input g_u with parameters set as *expanded_nodes*=2 and *depth*=2, a well-designed prompt strategy can enhance model performance. However, the ranking operation does not necessarily improve the performance of all LLMs; they have inconsistent preferences for different prompt strategies. For instance, ChatGPT favors the prompt method incorporating relevance scores, while the Vicuna series prefers ranking strategies. This discrepancy may be attributed to variations in the training data and inherent structure of the respective models. This finding suggests that when designing a prompt method, the universal applicability of a given prompt strategy across multiple LLMs should be considered.

Datasets	ChatGPT				Vicuna 7B				Vicuna 13B			
	g_u	Grouping	Ranking	Scoring	g_u	Grouping	Ranking	Scoring	g_u	Grouping	Ranking	Scoring
QALD-7	82.56	84.11	84.11	84.11	65.70	64.84	66.54	53.64	77.00	75.52	77.81	72.40
LC-QuAD 2.0	50.60	48.71	50.01	52.48	32.84	33.49	35.72	26.14	42.98	45.10	45.13	42.57
KQAPro	54.03	50.25	52.29	54.03	27.63	27.74	31.32	24.92	32.75	36.05	37.64	35.12

Table 4: Distinct models exhibit unique preferences towards various prompting methods.

5. Evaluating LLM’s Robustness to Noisy or Incomplete KG

In real-world scenarios, the subgraphs provided to LLMs often contain noise or are incomplete. Therefore, understanding the robustness of LLMs when handling noisy or incomplete subgraphs is crucial. In this section, we gradually corrupt the subgraphs to observe the performance of LLMs in linearized triple input formats.

5.1. The Generation of Noisy or Incomplete Subgraphs

As illustrated in Figure 6, we systematically corrupt the subgraph in two approaches: 1) Nodes are proportionally replaced with random irrelevant KG nodes. 2) Nodes are proportionally deleted randomly.

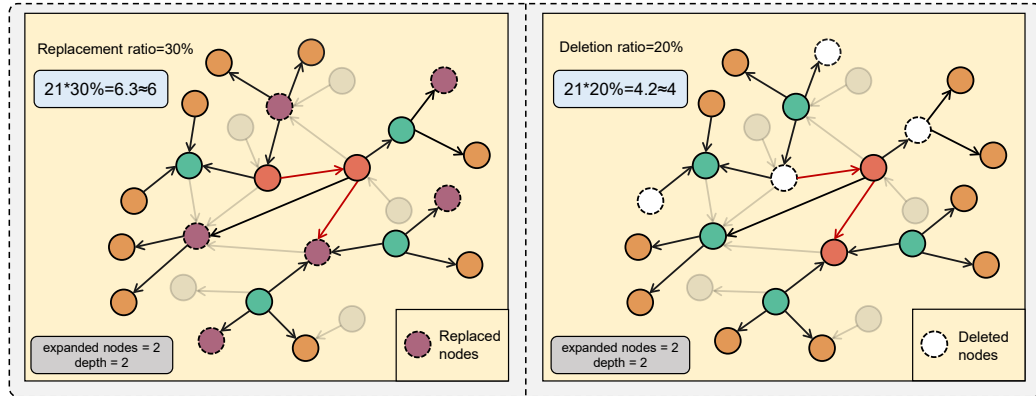


Figure 6: We calculate the number of nodes for random replacement and deletion based on the ratio. The nodes used for replacement are randomly sampled from unrelated nodes in the KG.

In the replacement operation, unrelated nodes replace parts of golden nodes to generate false fact information, denoted as a noisy subgraph. The deletion operation simulates an incomplete KG scenario. We randomly replace and delete nodes according to specified percentages. The replacement and deletion ratio ranges from 10% to 90%.

5.2. Datasets and Metric

Based on all datasets from the *Triple-to-Text* experiments, we utilize subgraphs with parameters set to *expanded_nodes=2* and *depth=2*, where nodes are randomly replaced or deleted. The evaluation metrics adhere to those used in the Section 3.1.1. We aim to observe the robustness of different LLMs when the proportion of replacements and deletions increases.

5.3. Results Analysis

As shown in Table 5, across three datasets, ChatGPT shows greater performance degradation on replacement and deletion compared to the smaller parameter models, Vicuna 7B and 13B. We have two preliminary findings: 1) The random replacement of nodes in KG has a more significant negative impact on the inference performance of LLM compared to the random deletion. *Incorrect facts are more likely to result in erroneous model outputs.* 2) Despite larger models demonstrating superior answering performance, they exhibit a more significant performance loss when facing random replacement and deletion of the subgraph. *There exists an inverse proportionality between a model’s robustness and its size.*

	QALD-7						LC-QuAD 2.0						KQAPro					
	ChatGPT		Vicuna 7b		Vicuna 13b		ChatGPT		Vicuna 7b		Vicuna 13b		ChatGPT		Vicuna 7b		Vicuna 13b	
Ratio	Delete	Replace	Delete	Replace	Delete	Replace	Delete	Replace	Delete	Replace	Delete	Replace	Delete	Replace	Delete	Replace	Delete	Replace
0%	82.56	82.56	65.70	65.70	77.00	77.00	50.60	50.60	32.84	32.84	42.98	42.98	54.03	54.03	27.63	27.63	32.75	32.75
10%	82.08	80.62	64.44	58.89	72.85	78.21	47.84	47.74	31.07	32.20	41.66	40.55	51.57	52.04	26.12	26.23	31.32	31.24
20%	82.08	80.62	70.00	58.65	75.51	81.16	46.23	46.02	31.83	29.59	39.43	39.43	49.66	49.91	26.09	26.17	30.96	31.07
30%	80.39	79.90	62.95	59.08	74.59	78.70	44.18	43.60	29.51	28.10	38.57	38.85	46.95	46.87	23.52	24.58	28.94	29.81
40%	80.39	79.90	55.27	57.68	77.00	73.84	42.49	42.41	27.83	27.47	36.12	36.83	43.82	44.10	23.60	23.94	27.93	28.19
50%	80.39	80.39	66.86	59.66	76.52	74.35	40.78	40.28	27.78	26.03	34.57	35.09	41.02	41.33	21.95	21.87	26.26	27.88
60%	80.39	80.39	63.24	61.30	75.31	78.21	36.25	35.99	24.50	22.56	33.50	32.94	37.28	37.08	21.37	19.16	25.53	25.36
70%	75.56	75.56	63.00	48.74	73.62	68.07	31.87	31.72	21.82	19.44	30.54	27.36	32.33	32.02	18.06	16.69	23.01	23.38
80%	65.85	63.91	50.00	47.25	65.17	66.14	26.28	25.82	18.54	16.50	25.56	23.79	27.21	27.46	14.57	13.49	20.41	20.25
90%	62.42	54.30	46.64	45.02	57.51	56.23	19.16	17.64	13.21	11.98	18.48	17.51	19.91	19.77	10.49	9.82	14.15	14.90
↓	20.14	28.26	19.06	20.68	19.49	20.77	31.44	32.96	19.63	20.86	24.50	25.47	34.12	34.26	17.14	17.81	18.60	17.85

Table 5: Randomly delete and replace nodes in the subgraph. ↓ quantifies the discrepancy between the model’s peak performance and its poorest performance.

6. Conclusion

In this study, we empirically study various input formats of KG injected into LLMs, yielding fundamental insights. For fact-intensive questions, linearized triples prompt LLMs more effectively than NL text. Additionally, we designed several simple prompt strategies to observe LLMs’ preference for the organization of triples, and the experiment results show that different series of LLMs do not

necessarily rely on order. We also evaluate the robustness of LLMs to noisy or incomplete subgraphs, discovering that larger models are more susceptible to these issues. The experiments suggest that when designing prompt strategies for larger LLMs, the impact of noise in the subgraph should also be considered. In summary, our findings offer valuable guidance for refining KG-related prompt strategies and underscore the importance of linearized triple knowledge in prompting LLMs.

Acknowledgments

This work is supported by the Natural Science Foundation of China (Grant No. U21A20488). We thank the Big Data Computing Center of Southeast University for providing the facility support on the numerical calculations in this paper.

References

- [1] S. Pan, L. Luo, Y. Wang, C. Chen, J. Wang, X. Wu, Unifying large language models and knowledge graphs: A roadmap, *IEEE Transactions on Knowledge and Data Engineering* 36 (2023) 3580–3599.
- [2] Y. Sun, S. Wang, Y. Li, S. Feng, X. Chen, H. Zhang, X. Tian, D. Zhu, H. Tian, H. Wu, Ernie: Enhanced representation through knowledge integration, *arXiv preprint arXiv:1904.09223* (2019).
- [3] W. Xiong, J. Du, W. Y. Wang, V. Stoyanov, Pretrained encyclopedia: Weakly supervised knowledge-pretrained language model, *arXiv preprint arXiv:1912.09637* (2019).
- [4] Y. Su, X. Han, Z. Zhang, Y. Lin, P. Li, Z. Liu, J. Zhou, M. Sun, Cokebert: Contextual knowledge selection and embedding towards enhanced pre-trained language models, *AI Open* 2 (2021) 127–134.
- [5] S. Arora, S. Wu, E. Liu, C. Ré, Metadata shaping: A simple approach for knowledge-enhanced language models, in: *Findings of the Association for Computational Linguistics: ACL 2022*, 2022, pp. 1733–1745. doi: 10.18653/v1/2022.findings-acl.137.
- [6] Q. Chen, F.-L. Li, G. Xu, M. Yan, J. Zhang, Y. Zhang, Dictbert: dictionary description knowledge enhanced language model pre-training via contrastive learning, *arXiv preprint arXiv:2208.00635* (2022).

- [7] F. Ufuk, The role and limitations of large language models such as chatgpt in clinical settings and medical journalism, *Radiology* 307 (3) (2023) e230276.
- [8] T. Sorensen, J. Robinson, C. M. Rytting, A. G. Shaw, K. J. Rogers, A. P. Delorey, M. Khalil, N. Fulda, D. Wingate, An information-theoretic approach to prompt engineering without ground truth labels, in: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, p. 819–862. doi:10.18653/v1/2022.acl-long.60.
- [9] J. White, Q. Fu, S. Hays, M. Sandborn, C. Olea, H. Gilbert, A. Elnashar, J. Spencer-Smith, D. C. Schmidt, A prompt pattern catalog to enhance prompt engineering with chatgpt, *arXiv preprint arXiv:2302.11382* (2023).
- [10] S. Li, Y. Gao, H. Jiang, Q. Yin, Z. Li, X. Yan, C. Zhang, B. Yin, Graph reasoning for question answering with triplet retrieval, in: *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9–14, 2023*, 2023, pp. 3366–3375. doi:10.18653/V1/2023.FINDINGS-ACL.208.
- [11] Y. Wen, Z. Wang, J. Sun, MindMap: Knowledge graph prompting sparks graph of thoughts in large language models, in: *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024, pp. 10370–10388. doi:10.18653/v1/2024.acl-long.558.
- [12] M. Sclar, Y. Choi, Y. Tsvetkov, A. Suhr, Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting, *arXiv preprint arXiv:2310.11324* (2023).
- [13] A. Voronov, L. Wolf, M. Ryabinin, Mind your format: Towards consistent evaluation of in-context learning improvements, in: *Findings of the Association for Computational Linguistics: ACL 2024, 2024*, pp. 6287–6310. doi:10.18653/v1/2024.findings-acl.375.
- [14] P. Zhan, Z. Xu, Q. Tan, J. Song, R. Xie, Unveiling the lexical sensitivity of llms: Combinatorial optimization for prompt enhancement, *arXiv preprint arXiv:2405.20701* (2024).
- [15] J. Baek, A. F. Aji, A. Saffari, Knowledge-augmented language model prompting for zero-shot knowledge graph question answering, *Proceedings of*

the First Workshop on Matching From Unstructured and Structured Data (MATCHING 2023) (2023).

- [16] P. Sen, S. Mavadia, A. Saffari, Knowledge graph-augmented language models for complex question answering, in: Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE), 2023, pp. 1–8.
- [17] K. Ma, H. Cheng, X. Liu, E. Nyberg, J. Gao, Open domain question answering with a unified knowledge interface, in: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2022, pp. 1605–1620. doi:10.18653/v1/2022.acl-long.113.
- [18] G. Xiong, J. Bao, W. Zhao, Y. Wu, X. He, Autoqgs: Auto-prompt for low-resource knowledge-based question generation from sparql, in: Proceedings of the 31st ACM International Conference on Information & Knowledge Management, 2022, pp. 2250–2259.
- [19] Y. Wu, N. Hu, G. Qi, S. Bi, J. Ren, A. Xie, W. Song, Retrieve-rewrite-answer: A kg-to-text enhanced llms framework for knowledge graph question answering, arXiv preprint arXiv:2309.11206 (2023).
- [20] T. Guo, Q. Yang, C. Wang, Y. Liu, P. Li, J. Tang, D. Li, Y. Wen, Knowledge-navigator: Leveraging large language models for enhanced reasoning over knowledge graph, *Complex & Intelligent Systems* 10 (5) (2024) 7063–7076.
- [21] P. Lai, F. Ye, Y. Fu, Z. Chen, Y. Wu, Y. Wang, V. Chang, Cognlg: Cognitive graph for kg-to-text generation, *Expert Systems* 41 (1) (2024) e13461.
- [22] P. Sahoo, A. K. Singh, S. Saha, V. Jain, S. Mondal, A. Chadha, A systematic survey of prompt engineering in large language models: Techniques and applications, arXiv preprint arXiv:2402.07927 (2024).
- [23] W. Liu, P. Zhou, Z. Zhao, Z. Wang, Q. Ju, H. Deng, P. Wang, K-bert: Enabling language representation with knowledge graph, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2020, pp. 2901–2908.
- [24] H. Zha, Z. Chen, X. Yan, Inductive relation prediction by bert, in: Proceedings of the AAAI conference on artificial intelligence, 2022, pp. 5923–5931.

- [25] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, P. Fung, Survey of hallucination in natural language generation, *ACM Computing Surveys* 55 (12) (2023) 1–38.
- [26] W. Chen, Large language models are few(1)-shot table reasoners, in: *Findings of the Association for Computational Linguistics: EACL 2023*, 2023, pp. 1120–1130. doi:10.18653/v1/2023.findings-eacl.83.
- [27] J. Jiang, K. Zhou, Z. Dong, K. Ye, W. X. Zhao, J.-R. Wen, StructGPT: A general framework for large language model to reason over structured data, in: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 9237–9251. doi:10.18653/v1/2023.emnlp-main.574.
- [28] H. Gao, L. Wu, P. Hu, F. Xu, Rdf-to-text generation with graph-augmented structural neural encoders, in: *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, 2021, pp. 3030–3036.
- [29] J. Wang, W. Huang, Q. Shi, H. Wang, M. Qiu, X. L. Li, M. Gao, Knowledge prompting in pre-trained language model for natural language understanding, in: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 3164–3177. doi:10.18653/v1/2022.emnlp-main.207.
- [30] T. Tang, J. Li, W. X. Zhao, J.-R. Wen, Mvp: Multi-task supervised pre-training for natural language generation, in: *Findings of the Association for Computational Linguistics: ACL 2023*, 2022, pp. 8758–8794. doi:10.18653/v1/2023.findings-acl.558.
- [31] D. Vrandečić, M. Krötzsch, Wikidata: a free collaborative knowledgebase, *Communications of the ACM* 57 (10) (2014) 78–85.
- [32] R. Usbeck, A.-C. N. Ngomo, B. Haarmann, A. Krithara, M. Röder, G. Napolitano, 7th open challenge on question answering over linked data (qald-7), in: *Semantic Web Challenges: 4th SemWebEval Challenge at ESWC 2017*, Portoroz, Slovenia, May 28-June 1, 2017, Revised Selected Papers, Springer, 2017, pp. 59–69.
- [33] M. Dubey, D. Banerjee, A. Abdelkawi, J. Lehmann, Lc-quad 2.0: A large dataset for complex question answering over wikidata and dbpedia, in: *The*

Semantic Web–ISWC 2019: 18th International Semantic Web Conference, Auckland, New Zealand, October 26–30, 2019, Proceedings, Part II 18, 2019, pp. 69–78.

- [34] S. Cao, J. Shi, L. Pan, L. Y. Nie, Y. Xiang, L. Hou, J. Li, B. He, H. Zhang, KQA pro: A dataset with explicit compositional programs for complex question answering over knowledge base, in: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2022, pp. 6101–6119.
- [35] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, et al., Judging llm-as-a-judge with mt-bench and chatbot arena, *Advances in Neural Information Processing Systems* 36 (2023) 46595–46623.
- [36] P. Rajpurkar, J. Zhang, K. Lopyrev, P. Liang, SQuAD: 100,000+ questions for machine comprehension of text, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, 2016, pp. 2383–2392. doi:10.18653/v1/D16-1264.
- [37] Y. Tan, D. Min, Y. Li, W. Li, N. Hu, Y. Chen, G. Qi, Can chatgpt replace traditional kbqa models? an in-depth analysis of the question answering performance of the gpt llm family, in: International Semantic Web Conference, Springer, 2023, pp. 348–367.
- [38] Y. Yao, D. Ye, P. Li, X. Han, Y. Lin, Z. Liu, Z. Liu, L. Huang, J. Zhou, M. Sun, Docred: A large-scale document-level relation extraction dataset, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 2019, pp. 764–777. doi:10.18653/v1/P19-1074.
- [39] D. Zhang, S. Mohan, M. Torkar, A. McCallum, A distant supervision corpus for extracting biomedical relationships between chemicals, diseases and genes, in: Proceedings of the Thirteenth Language Resources and Evaluation Conference, 2022, pp. 1073–1082.
- [40] Y. Zhang, H. Dai, Z. Kozareva, A. Smola, L. Song, Variational reasoning for question answering with knowledge graph, in: Proceedings of the AAAI conference on artificial intelligence, 2018, pp. 6069 – 6076.
- [41] L. Wang, L. Li, D. Dai, D. Chen, H. Zhou, F. Meng, J. Zhou, X. Sun, Label words are anchors: An information flow perspective for understanding

- in-context learning, in: Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, 2023, pp. 9840–9855. doi:10.18653/v1/2023.emnlp-main.609.
- [42] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 2019, pp. 4171–4186. doi:10.18653/v1/N19-1423.
- [43] L. Logeswaran, M.-W. Chang, K. Lee, K. Toutanova, J. Devlin, H. Lee, Zero-shot entity linking by reading entity descriptions, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 2019, pp. 3449–3460.
- [44] S. Humeau, K. Shuster, M.-A. Lachaux, J. Weston, Poly-encoders: Transformer architectures and pre-training strategies for fast and accurate multi-sentence scoring, arXiv preprint arXiv:1905.01969 (2019).