

Thinking with Knowledge Graphs: Enhancing LLM Reasoning Through Structured Data

Xue Wu
Yahoo Research
Mountain View, California, USA
xuewu@yahooinc.com

Kostas Tsioutsoulis
Facts.ai
Saratoga, California, USA
kostas@facts.ai

ABSTRACT

Large Language Models (LLMs) have demonstrated remarkable capabilities in natural language understanding and generation. However, they often struggle with complex reasoning tasks and are prone to hallucination. Recent research has shown promising results in leveraging knowledge graphs (KGs) to enhance LLM performance. KGs provide a structured representation of entities and their relationships, offering a rich source of information that can enhance the reasoning capabilities of LLMs. For this work, we have developed different techniques that tightly integrate KG structures and semantics into LLM representations. Our results show that we are able to significantly improve the performance of LLMs in complex reasoning scenarios, and ground the reasoning process with KGs. We are the first to represent KGs with programming language and fine-tune pretrained LLMs with KGs. This integration facilitates more accurate and interpretable reasoning processes, paving the way for more advanced reasoning capabilities of LLMs.

1 INTRODUCTION

Large Language Models (LLMs) have achieved state of art performance in many Natural Language Processing (NLP) tasks ([7], [11]). They have been successfully applied in a wide range of applications, from question answering to summarization and machine translation. However, due to limitations in the training data and training process, they suffer from hallucination [18], where generated text is nonfactual, nonsensical, or incoherent. This issue becomes particularly prevalent when dealing with tasks that require intricate or complex reasoning. To address hallucination, researchers have used methods such as prompting ([31], [39], [30], [14], [25], [35], [34]), retrieval augmented generation (RAG [20]), and fine-tuning. These approaches often leverage external sources of information, which can come from the internet, third-party applications, databases, or knowledge graphs.

Knowledge Graphs (KGs) are structured representations of real world entities and the relationships among them, offering a rich source of factual information ([21] [4]). By grounding the reasoning processes of LLMs with KGs, we can enhance the factual accuracy of the generated text and reduce hallucinations. Researchers have explored various approaches to integrate KGs with LLMs, each with its own advantages and limitations. One approach involves using Graph Neural Networks (GNNs) to encode KGs into embeddings that capture the structural and semantic information of the graphs ([10], [36]). These embeddings then serve as soft prompts to LLMs, guiding the generation process with knowledge from the KGs. However, this method requires very careful design and tuning, as it needs to align the representations learned by GNNs with the token-based processing of LLMs. Moreover, since such integration

is specific per graph, it may require significant re-engineering for different tasks or graph types.

Another method uses semantic parsing to convert natural language queries into structured query languages like SPARQL ([32], [6], [26], [8]). In this approach, the LLM generates a SPARQL query based on the input prompt, which is then executed against the KG to retrieve relevant information. This method effectively uses the KG as an external knowledge base, and treats the KG and the LLM as separate components. As a result, the reasoning process is not fully integrated into the LLM, potentially limiting the model's ability to perform complex reasoning during text generation.

Alternatively, researchers encode entities and relationships of KGs as natural language text ([22], [12]). They either incorporate them into the LLM's input context or fine-tune the LLM with these text representations. This approach leverages the LLM's natural language understanding ability to reason over the text-encoded knowledge. However, representing structured data as unstructured text poses challenges. Capturing the nuances of entities and relationships in natural language requires careful design to avoid ambiguities and ensure that the structural information is preserved, which can be non-trivial for complex graphs.

While there are different ways to leverage KGs to enhance LLMs' performance, there are two aspects to using KGs for improving LLMs. One is grounding the LLMs with trustworthy information, and the other is providing them with examples of relations from which they can generalize. It is this second part that motivated our work. In this work, we propose an approach that represents knowledge graphs with programming language code. Programming languages are inherently structured and are designed to represent complex data and relationships efficiently. This allows for an accurate encoding of entity relationships that preserves the internal structure of the graph. More importantly, programming code is part of the pre-training data for many LLMs, meaning that the models are already equipped to parse and understand programming syntax and semantics. This reduces the need for additional specialized training to interpret the KG representations. By leveraging the structured representation of KGs in programming languages, LLMs can perform more sophisticated reasoning over the data. We investigate different methods of integrating entity relationships into LLMs. Our experimental results demonstrate that programming language (Python) representations of KGs outperform traditional natural language representations and structured JSON representations in complex reasoning tasks, leading to more accurate and reliable outputs.

The main contributions of this paper are:

- We introduce a novel representation of knowledge graphs with programming language. It facilitates the seamless integration of structured knowledge into the language modeling process.
- By tightly integrating knowledge graphs into LLMs, our approach improves the reasoning accuracy of LLMs on complex tasks and effectively grounds the reasoning process, reducing the chance for hallucinations.

2 RELATED WORK

There have been several attempts to apply LLMs to graph reasoning tasks. Wang et al. [28], Guo et al. [15], and Ye et al. [37] employ the Graph2Text strategy of converting graph data into textual descriptions. However, these textual descriptions can result in very large contexts, and algorithms such as shortest path computations and bipartite graph matching require calculations across the entire context, making the task highly challenging. Chai et al. [9] have introduced GraphLLM, which combines three steps, namely node understanding, graph structure understanding, and graph-enhanced prefix tuning. Zhu et al. [40], and Wang et al. [29] proposed different methods for instruction fine-tuning LLMs to improve the performance of common graph tasks.

While the above works address general graph problems, other research has focused specifically on combining KGs with LLMs. One such approach is to use the LLM as an encoder to transform text-based entity nodes and relations, and then fuse the LLM and GNN-derived representations. Applications of this approach range from product recommendation (Choudhary et al. [10]) to biomedical question-answering (Yasunaga et al. [36]). Luo & Pan [22] have proposed a reasoning on Graph (RoG) method that comprises two modules: a planning module and a retrieval-reasoning module. The planning module mines the KG and generates faithful relation paths for answering the question. The retrieval-reasoning module combines a breadth-first search and a probabilistic optimization over all paths. Dernbach et al. [12] developed a neighborhood partitioning and encoding scheme to accommodate real-world graph properties. Their encoding scheme transforms graph relations into alternate text representations, which in turn are used to fine-tune the LLM.

Edge et al. [13] have built an end-to-end system, called GraphRAG, that starts with a set of source documents and iteratively applies an LLM to extract the entities in each document. Next, entities are connected via the relationships extracted, and a knowledge graph is created. The knowledge graph is split into communities, and each community is summarized. These summaries are subsequently used by the LLM via RAG to help answer questions submitted to the system.

Nie et al. [24] provide a code-style in-context learning (ICL) method for knowledge base question answering (KBQA). They design seven meta-functions written in Python that cover the atomic operations used for querying databases. By using few-shot learning, they improve LLMs’ ability to query knowledge bases effectively.

There is also ongoing research work ([23], [38], [2]) that studies the impact of mixing programming code in pre-training or instruction fine-tuning datasets on the performance of LLMs. Even though the programming code used by the research is generic, the results

show promising improvements on various tasks, including reasoning tasks. There is another line of work ([19], [39]) that studies how to improve LLM Chain of Thought (CoT) reasoning ability through fine-tuning. However, creating good datasets for fine-tuning CoT is labor intensive, and the reasoning steps are described in natural language text, which can introduce ambiguity.

Yang et al. [33] studied how LLMs do multi-hop reasoning and found that LLMs perform latent multi-hop reasoning in certain relation types, but often struggle with later hops. Biran et al. [3] observed that later hops are resolved in the model’s deeper layers, and thus the LLM may no longer encode the necessary knowledge for answering the question. They propose a back-patching approach, essentially feeding the hidden representations from later layers back into earlier ones.

In this work, we continue the prior research work on multi-hop queries, showing that we can significantly improve the performance of LLMs by integrating KG structures and semantics into LLM representations. Similarly to [24] we also experiment with code-style representations. However, in our case, we represent the entity-relations -not the atomic operations- in Python, and our goal is to improve the LLMs’ ability to answer questions directly, not its ability to query the knowledge graph.

The rest of the paper is structured as follows. In Section 3 we describe the methodologies we followed to prompt or fine-tune the LLM. In Section 4 we describe the experimental design and results. Finally, we conclude the paper in Section 5.

3 METHODOLOGY

Our work focuses on studying the entity relationships representation of KG for grounded LLM reasoning.

3.1 Knowledge Graph Definition

Let $G = \{E, R, T\}$ denote a knowledge graph, where E is the set of entities, R is the set of relationships, $T \subseteq E \times R \times E$ is the set of triplets that are edges in the knowledge graph. A triplet $(e_i, r_i, e_{i+1}) \in T$ if there is a directed edge between entity e_i ($e_i \in E$) and entity e_{i+1} ($e_{i+1} \in E$) through relationship r_i ($r_i \in R$). A triplet also corresponds to a complete one hop reasoning process. A two hop compositional reasoning process can be represented as $((e_i, r_i, e_{i+1}), (e_{i+1}, r_{i+1}, e_{i+2}))$, where $(e_i, r_i, e_{i+1}) \in T$ and $(e_{i+1}, r_{i+1}, e_{i+2}) \in T$. Since e_{i+1} exists in both triplets, it is the bridge entity. Similarly, a three hop compositional reasoning process can be represented as $((e_i, r_i, e_{i+1}), (e_{i+1}, r_{i+1}, e_{i+2}), (e_{i+2}, r_{i+2}, e_{i+3}))$, where $(e_i, r_i, e_{i+1}) \in T$, $(e_{i+1}, r_{i+1}, e_{i+2}) \in T$ and $(e_{i+2}, r_{i+2}, e_{i+3}) \in T$.

3.2 Knowledge Graph Representation for LLM Reasoning

To improve LLM multihop reasoning with knowledge graphs, we represent knowledge graphs in ways that are more compatible with LLM prompting and fine-tuning. When given a complex reasoning prompt, LLMs can detect entities and relationships, then implicitly infer the key entities and facts by following logical reasoning steps grounded by knowledge graphs. For instance, given the prompt: “Who is the spouse of the composer of ‘It Goes Like It Goes?’”, LLMs can follow the reasoning steps: “The composer of ‘It Goes Like It Goes’ is David Shire”, and “The spouse of David Shire is Didi

Conn", and infer that the correct answer is "Didi Conn". Different representations of knowledge graphs affect how effectively LLMs can perform such logical reasoning.

The most natural way is to use natural language to describe the triplets in knowledge graphs. Figure 1 shows the natural language representation for two-hop reasoning of triplets $((e_1, r_1, e_2), (e_2, r_2, e_3))$, where e_2 is the bridge entity. For instance, the triplets (('It Goes Like It Goes', 'composer', 'David Shire'), ('David Shire', 'spouse', 'Didi Conn')) can be represented as (The composer of 'It Goes Like It Goes' is David Shire, The spouse of David Shire is Didi Conn). And the two hop reasoning can be represented as "The spouse of the composer of 'It Goes Like It Goes' is Didi Conn". As LLMs understand natural languages very well, the natural language representation is the most straightforward way for either prompting or fine-tuning LLMs with knowledge graphs.

JSON (JavaScript Object Notation) is a lightweight data interchange format [5]. It is designed to store data in universal data structures such as dictionary and list that are supported by most programming languages. JSON is a pure data only format and can be used to store structured data from knowledge graphs. Figure 2 shows the JSON representation of knowledge graph triplets (e_1, r_1, e_2) and (e_2, r_2, e_3) . The entities e_1 and e_2 are the keys and the relationship/entity pairs $r_1 : e_2$ and $r_2 : e_3$ are the values. However, since JSON is designed to store data only, it is difficult to represent multi-hop inference process with the JSON format. Alternatively, a comment or description field (with "comment" as the key and natural language description of the multi-hop reasoning as the value) can be added to the JSON representation. But this is not a recommended practice in general.

Programming language code is another major data source for LLM pre-training and fine-tuning. Knowledge graphs can be represented using various data structures supported by major programming languages such as Python. The triples can be represented either as a static dictionary or added dynamically to the predefined data structures as part of the code. Figure 3 shows Python representation of knowledge graph triplets with a static dictionary data structure, and the two-hop inference process that is based on the stored triplets. As shown in figure 3, relationships and entities are stored in dictionary data structure "relationships" with relationships r_1 and r_2 as the keys and entity pairs as the values. The inference process is just the process of retrieving the values based on the key values of the dictionary. Figure 4 shows Python representation of knowledge graph triplets with predefined Python class "KnowledgeBase", and an iterative multi-hop inference function 'infer' that supports inference of an arbitrary number of hops. As shown by figure 4, the main data structure of "KnowledgeBase" is a dictionary "self.facts", and entities and relationships are added to the dictionary with function "add_fact". The "infer" function accepts any number of relationships with parameter "*relations" and does the corresponding multi-hop reasoning. We designed the dynamic self-defined data structure based Python representation because it can be easily generalized to support multi-hop reasoning based on subgraphs of KGs.

The example mentioned in the previous paragraphs can be easily represented as dictionaries in Python. Figure 5 shows four different representations for the same example. Using programming languages, knowledge graphs can be represented in a more controlled

and unambiguous way; corner cases can be easily checked. However, the representation for the same triplet can be more verbose.

r1 of e1 is e2, r2 of e2 is e3. r2 of r1 of e1 is e3.

Figure 1: Natural Language Representation of KG with Static Relationships

```
{
  "{r1}": "{e1}": "{e2}",
  "{r2}": "{e2}": "{e3}"
}
```

Figure 2: JSON Representation of KG with Static Relationships

```
# Step 1. Define relationships and entities
relationships = {
  'r1': {'e1': 'e2'},
  'r2': {'e2': 'e3'}
}

e1 = 'e1'
r1 = 'r1'
r2 = 'r2'

# Step 2. Perform inference
e2 = relationships[r1][e1]
e3 = relationships[r2][e2]
```

Figure 3: Python Representation of KG with Static Relationships

```
# Step 1. Define relationships with knowledge base
class KnowledgeBase:
    def __init__(self):
        self.facts = {}

    def add_fact(self, entity1, relation, entity2):
        self.facts[(entity1, relation)] = entity2

    def infer(self, entity, *relations):
        current_entity = entity
        for relation in relations:
            key = (current_entity, relation)
            if key in self.facts:
                current_entity = self.facts[key]
            else:
                # If the path does not exist, return None.
                return None
        return current_entity

kb = KnowledgeBase()
kb.add_fact(e1, r1, e2)
kb.add_fact(e2, r2, e3)

# Step 2. Perform inference.
result = kb.infer(e1, r1, r2)
```

Figure 4: Python Representation of KG with Dynamic Relationships

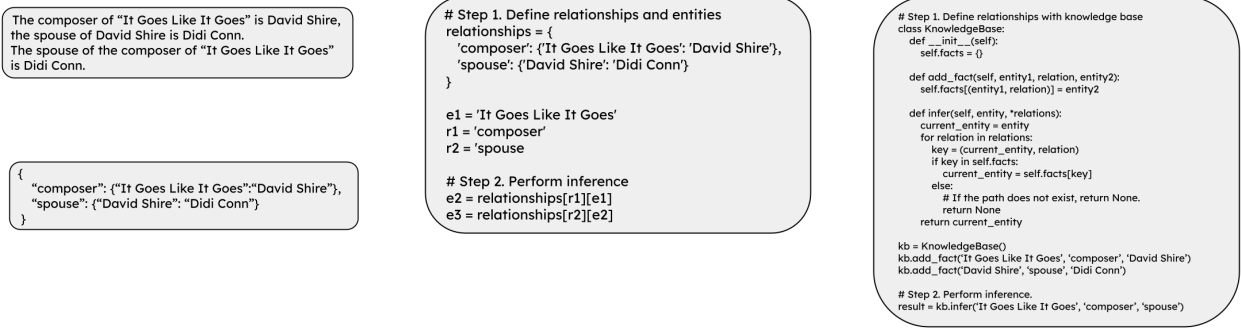


Figure 5: Examples of Different Representation of KG

	Train	Test	Train and Test Intersection
Number of Hops	2	2	2
Dataset Size	10,262	10,255	0
Bridge Entities (e_2)	324	880	0
Relations (r_1, r_2)	191	207	167
No. of row with (r_1, r_2)	10,262	10,255	10,130

Table 1: Dataset 1 Train and Test Data Selection

	Train	Test	Train and Test Intersection
Number of Hops	2	2	2
Dataset Size	10,733	5,236	0
Bridge Entities (e_2)	4,170	3,650	782
Relations (r_1, r_2)	80	110	66
No. of row with (r_1, r_2)	10,733	5,236	1,907

Table 2: Dataset 2 Train and Test Data Selection

4 EXPERIMENTS

We designed experiments to study how different representations of entity relationships in KGs affect the reasoning performance of LLMs across two different datasets.

4.1 Datasets

4.1.1 Dataset 1. For the first dataset, we use the same dataset used by the “Hopping Too Late” paper ([3]). The dataset includes two-hop relationships extracted from publicly available knowledge base Wikidata ([27]). For the experiments, we split the dataset into eight equal-sized partitions based on the bridge entity e_2 in a round-robin fashion, so that any unique e_2 exists only in one partition. We partition the dataset in this way so that the LLMs only learn the relationships and logical reasoning process rather than memorizing the entities. To avoid overrepresentation of the most popular e_2 in the training or testing dataset, we choose partition 2 as our training dataset and partition 4 as the testing dataset. The details of the dataset are listed in Table 1. For this dataset, there is no overlap of bridge entities between training and testing dataset. The overlap of relationship pairs (r_1, r_2) is about 99%.

4.1.2 Dataset 2. For the second dataset, we use the dataset created in paper ([16]). This dataset also includes two-hop relationships extracted from publicly available knowledge base Wikidata ([27]). Although both Dataset 1 and Dataset 2 are derived from the same knowledge base, the extracted entities and relationships for Dataset 2 are different from those in Dataset 1. For the training dataset, we select only compositional relationships from this dataset and limit the number of instances per relationship pair (r_1, r_2) to no more than 500 to avoid over representation of particular relationships in

the training data. We use the development dataset with compositional relationships for testing purpose. The details of the dataset are listed in Table 2. We choose compositional relationships so the type of relationships are consistent with those in Dataset 1 and it is easier to compare the reasoning performance across datasets. We didn’t further restrict the entities and relationships based on the overlap between training and testing dataset to respect the design of the dataset by paper ([16]).

4.1.3 Dataset 3. This dataset is an extension of Dataset 1. We extended two-hop relationships $((e_1, r_1, e_2), (e_2, r_2, e_3))$ by adding a third hop (e_3, r_3, e_4) , resulting in three-hop relationships $((e_1, r_1, e_2), (e_2, r_2, e_3), (e_3, r_3, e_4))$, while keeping the entities e_3 as a subset of the entities e_3 in Dataset 1. This dataset was created to test whether models fine-tuned for two-hop reasoning can generalize to improve three-hop reasoning performance as well. The details of the dataset are listed in Table 3. The overlap of bridge entities between training and testing dataset is minimal for this dataset. There is high percentage of relationship overlap between training and testing dataset.

The final training data for fine-tuning the LLMs include both one-hop prompts and responses based on the first hop from the datasets, and two-hop prompts and responses based on two-hop information from the datasets.

4.2 Large Language Models

To evaluate the performance of LLM reasoning, we chose the latest released open source models by Meta: Llama-3.1-8B-Instruct and Llama-3.1-70B-Instruct ([1]). The Llama-3.1 model family demonstrates stronger reasoning abilities compared to other open-source

	Dataset 1 Train	Test	Intersection with Dataset 1 Train
Number of Hops	2	3	0
Dataset Size	10,262	1,007	0
Bridge Entities e_2	324	261	28
Train e_2 vs Test e_3	324	233	15
Relations (r_1, r_2)	191	81	70
No. of row with (r_1, r_2)	10,262	1,007	973
Relations (r_2, r_3)	*191	97	92

Note *: this number is for (r_1, r_2) in Dataset 1 Train

Table 3: Dataset 3 Basic Statistics

LLMs. For our experiments, we fine-tuned the smaller model (Llama-3.1-8B-Instruct) with different entity relationship representations of multi-hop reasoning. During fine-tuning, we optimize the language model objective, which predicts the probability distribution of each token based on its previous tokens. We adopt LoRA ([17]) to perform parameter-efficient fine-tuning. During inference, we used greedy decoding for reproducibility.

4.3 Experiment Environment

We use the large language model checkpoints from Hugging Face Transformers¹ in all the experiments. Experiments were conducted using one Nvidia A100 processor with 40GB of RAM per GPU. When fine-tuning the LLMs, we ran the training process for one epoch across all the experiments. We didn’t fine-tune the larger model (Llama-3.1-70B-Instruct) because of the computation resource constraint.

4.4 Experiment Design

To test how well each representation of KGs helps with LLM reasoning, we designed the following three types of experiments.

4.4.1 LLM Reasoning without Context. To test the multi-hop reasoning ability of LLMs, we compared model performance under three conditions:

- (1) Zero shot prompting of the LLMs to predict e_3 . Given (e_1, r_1, r_2) , we prompt the LLMs to predict e_3 . This serves as our baseline experiment.
- (2) One shot prompting of the LLMs with reasoning example. We have conducted three experiments, one for each representation format of KGs.
- (3) Zero shot prompting of fine-tuned LLMs. We have conducted three experiments, one for each representation format of KGs.

4.4.2 LLM Reasoning with Context. Since Retrieval Augmented Generation (RAG) is one of the most popular applications of LLMs, we designed experiments to compare the performance of LLMs when given input context. This set of experiments are designed to test how well LLMs perform for potential RAG applications.

4.4.3 Model Generalization to Longer Reasoning Path. We fine-tuned the LLMs using only one-hop and two-hop triplets. For this set of experiments, we tested how well the fine-tuned models can generalize to reasoning over longer paths.

¹<https://huggingface.co/docs/transformers/index>

4.5 Prompt Design

The design of the prompts for all the experiments is presented in Table 4. Examples are shown for two-hop reasoning prompts. Because of the paper length limitation, the natural language explanation, JSON data structure and Python code snippets for one-shot prompting are shown in Appendix A. Three-hop reasoning prompts follow the same logic and use the same prompt format.

4.6 Metrics

The main metric for measuring the performance of LLM reasoning is the accuracy of multi-hop reasoning conditioned on the correctness of each individual hop, denoted as $r = p(h|h_1, h_2, \dots, h_n)$, where h is the correctness of the multi-hop reasoning, and h_i is the correctness of the i th hop reasoning. We use this metric instead of the overall accuracy of the results because the final reasoning output is affected by multiple factors, including whether the LLM has the knowledge of each individual entity and whether it can infer each hop correctly. Our main metric measures the latent multi-hop reasoning performance. We also provide the overall accuracy of the results as a reference for potential applications.

4.7 Experimental Results

We compare the different approaches of grounding LLM reasoning using different knowledge graph representations by answering the following research questions.

4.7.1 RQ1: Do different representations of entity relationships affect LLM multi-hop reasoning? Can we fine-tune the LLM with proper entity relationship representations to improve its reasoning capabilities? Since LLMs have already demonstrated abilities in latent multi-hop reasoning ([33]), we designed experiments to compare how different approaches and different representations of multi-hop entity relationships can further improve LLMs’ multi-hop reasoning ability.

Impact of Entity Relationship Representations for LLM Prompting. Table 5 and Table 6 present the performance of Llama-3.1-8B-Instruct with different representations across different datasets. It is obvious from the one-shot prompting results that Python representations of entity relationships outperform both the plain natural language text representation and the JSON representation. And both Python and natural language representations for one shot prompting perform better than zero shot prompting. The two-hop reasoning performance of the one-shot LLM with dynamic Python representation is approximately 78% higher than that of zero-shot prompting of the LLM on the Dataset 1 (Table 1). Similarly, the performance of the static Python representation for one-shot prompting of the LLM is about 60% higher than that of zero-shot prompting for Dataset 2 (Table 2). However, the performance of one-shot prompting with JSON representation is worse than zero-shot prompting of LLM. We hypothesize that the structured JSON representation is not native to the LLM reasoning process.

Impact of Entity Relationship Representation for LLM Fine-Tuning. Table 5 and Table 6 also show the performance of the fine-tuned Llama-3.1-8B-Instruct model with different representations on different datasets. In this case, the LLMs fine-tuned with Python representations perform better than the LLMs fine-tuned with either

Dataset	Dataset1	Dataset2
Multihop Question	r_2 of r_1 of e_1 is	What is r_2 of r_1 of e_1 ?
Zero Shot Prompt	Given the incomplete statement: r_2 of r_1 of e_1 is $_{-}$, provide answer and generate explanation for completing the statement	Given the question: What is r_2 of r_1 of e_1 ? generate explanation and provide answer to the question
One Shot Prompt for text representation	Given the incomplete statement: r_2 of r_1 of e_1 is $_{-}$, {"Answer": " e_3 ", "Explanation": " r_1 of e_1 is e_2 . r_2 of e_2 is e_3 . r_2 of r_1 of e_1 is e_3 "} Given the incomplete statement: {statement} $_{-}$, provide answer and generate explanation for completing the statement	Given the question: What is r_2 of r_1 of e_1 ? {"Answer": " e_3 ", "Explanation": " r_1 of e_1 is e_2 . r_2 of e_2 is e_3 . r_2 of r_1 of e_1 is e_3 "} Given the question: {question} ? generate explanation and provide answer to the question generate explanation and provide answer to the question
Or for JSON representation	Given the incomplete statement: r_2 of r_1 of e_1 is $_{-}$, {"Answer": " e_3 ", "JSON Structure": "{JSON data}" Given the incomplete statement: {statement} $_{-}$, provide answer and generate JSON structure for completing the statement	Given the question: What is r_2 of r_1 of e_1 ? {"Answer": " e_3 ", "JSON structure": "{JSON data}" Given the question: {question} ? generate JSON structure and provide answer to the question
Or for Python representation	Given the incomplete statement: r_2 of r_1 of e_1 is $_{-}$, {"Answer": " e_3 ", "Python code snippet": "{Python code}" Given the incomplete statement: {statement} $_{-}$, provide answer and generate python code for completing the statement	Given the question: What is r_2 of r_1 of e_1 ? {"Answer": " e_3 ", "Python code snippet": "{Python code}" Given the question: {question} ? generate python code and provide answer to the question
Prompt with Context	Given context: {context} and the uncompleted statement: {statement} $_{-}$, provide answer and generate explanation for completing the statement	Given context: {context} and the question: {question} generate explanation and provide answer to the question

Table 4: Prompt Design

JSON representation or plain natural language text representations. The performance of the LLM fine-tuned with JSON representation is only slightly worse than that of the LLM fine-tuned with natural language representation. For comparison, we also provide the performance numbers for zero shot prompting of Llama-3.1-70B-Instruct. All LLMs that were fine-tuned with any entity relationship representation outperform the larger model for latent multi-hop reasoning.

4.7.2 RQ2: Can fine-tuned LLMs generalize their reasoning ability to more hops that are beyond the training data ? Since we fine-tuned LLMs using only one-hop and two-hop reasoning data, it is important to study whether such a fine-tuning process will improve the LLMs’ reasoning performance on more hops. We created a three-hop dataset (as shown in Table 3) to measure the performance of the fine-tuned LLMs. As shown in Table 7, the three-hop reasoning performance of all fine-tuned LLMs has improved across all entity relationship representations compared with the baseline LLM without fine-tuning. Furthermore, LLMs that were fine-tuned with the Python representation outperform those fine-tuned with either the plain natural language representation or the JSON representation. As the relationship (r_2, r_3) has significant overlap with (r_1, r_2) in the training data (as shown in Table 3), the relative performance of fine-tuned models is consistent with what is shown in Table 5.

4.7.3 RQ3: How much can LLM in-context learning help or benefit from multi-hop reasoning? Retrieval Augmented Generation (RAG) is one of the main applications of LLMs. However, even when supplied with correctly retrieved information, LLMs do not always generate the correct answer, particularly when multi-hop reasoning is required. To address this issue, we designed experiments to study the performance of fine-tuned LLMs when given an input context.

In these experiments, we performed one-shot prompting of the LLMs, using the prompts listed in Table 4. The results are shown in Table 8. For Dataset 1, since no context is provided, we use the 1st hop and 2nd hop inferences as context and measure whether the model can infer the correct answers. For Dataset 2, we use the given context for the questions. And again, the LLMs that were fine-tuned with different entity relationship representations outperform the baseline model without fine-tuning. Notably, the LLM fine-tuned with dynamic Python representation greatly outperforms the baseline model; it even surpasses the performance of the much larger baseline model (Llama-3.1-70B-Instruct).

4.8 Discussion

We designed a series of experiments to study how to seamlessly integrate entity relationships into LLMs to improve their multi-hop reasoning ability, and the impact of different entity relationship representations on the performance of LLMs. Our experimental results show that it is possible to integrate entity relationships into LLMs and ground the LLM multi-hop reasoning process with knowledge graphs. While all forms of the proposed entity relationship representations help improve LLM reasoning performance, they affect LLM performance differently. The natural language representation is straightforward and is the major form of pre-training data for LLMs. The JSON representation is best suited for storing structured data. However, it can be difficult for LLM to directly integrate pure structured data. The Python representations store both structured data and the inference process, providing a more controlled and unambiguous way of guiding LLMs through the reasoning process. This helps LLMs achieve better reasoning performance in all the

Model	KG Representation for fine tuning	KG Representation for Prompt	Prompt	Accuracy	1st hop correct 2nd hop correct final incorrect	1st hop correct 2nd hop correct final correct	1st hop correct 2nd hop correct final accuracy
Llama-3.1-8B	NA	NA	Zero Shot	16.3%	989	604	38.2%
Llama-3.1-70B	NA	NA	Zero Shot	33.9%	615	1,979	76.3%
Llama-3.1-8B	NA	Natural Language	One Shot	16.9%	965	776	44.6%
Llama-3.1-8B	NA	JSON	One Shot	7.78%	995	352	26.1%
Llama-3.1-8B	NA	Python Static	One Shot	17.5%	408	609	59.9%
Llama-3.1-8B	NA	Python Dynamic	One Shot	19.1%	341	772	67.9%
Llama-3.1-8B tuned	Natural Language	NA	Zero Shot	25.6%	145	1,005	87.4%
Llama-3.1-8B tuned	JSON	NA	Zero Shot	20.8%	82	569	87.4%
Llama-3.1-8B tuned	Python Static	NA	Zero Shot	26.2%	102	1,093	91.5%
Llama-3.1-8B tuned	Python Dynamic	NA	Zero Shot	26.5%	119	882	88.1%

Table 5: LLM Two Hop Reasoning for Dataset 1

Model	KG Representation for fine tuning	KG Representation for Prompt	Prompt	Accuracy	1st hop correct 2nd hop correct final incorrect	1st hop correct 2nd hop correct final correct	1st hop correct 2nd hop correct final accuracy
Llama-3.1-8B	NA	NA	Zero Shot	1.85%	51	14	21.5%
Llama-3.1-70B	NA	NA	Zero Shot	10.7%	164	195	54.3%
Llama-3.1-8B	NA	Natural Language	One Shot	3.99%	82	25	23.4%
Llama-3.1-8B	NA	JSON	One Shot	1.80%	19	4	17.4%
Llama-3.1-8B	NA	Python Static	One Shot	4.98%	50	29	36.7%
Llama-3.1-8B	NA	Python Dynamic	One Shot	4.12%	35	12	25.5%
Llama-3.1-8B tuned	Natural Language	NA	Zero Shot	11.0%	97	164	62.8%
Llama-3.1-8B tuned	JSON	NA	Zero Shot	10.0%	61	157	72.0%
Llama-3.1-8B tuned	Python Static	NA	Zero Shot	12.1%	44	191	81.3%
Llama-3.1-8B tuned	Python Dynamic	NA	Zero Shot	12.3%	38	203	84.2%

Table 6: LLM Two Hop Reasoning for Dataset 2

Model	KG Representation for tuning	Accuracy	% of Correct conditioned on 1st & 2nd hop correct	% of Correct conditioned on 2nd & 3rd hop correct	% of Correct conditioned on all three hops correct
Llama-3.1-8B	NA	20.4%	35.2%	37.1%	50.0%
Llama-3.1-70B	NA	42.6%	64.6%	61.8%	80.4%
Llama-3.1-8B tuned	Natural Language	31.0%	44.0%	54.1%	65.0%
Llama-3.1-8B tuned	JSON	23.3%	42.5%	46.1%	63.9%
Llama-3.1-8B tuned	Python Static	30.9%	47.7%	59.3%	70.6%
Llama-3.1-8B tuned	Python Dynamic	38.0%	51.6%	57.4%	67.0%

Table 7: LLM Three Hop Reasoning for Dataset 3 with Zero Shot Prompting

Model	KG Representation for tuning	Dataset 1 Accuracy given 1st & 2nd hop as context	Dataset 2 Accuracy given context
Llama-3.1-8B	NA	88.6%	10.9%
Llama-3.1-70B	NA	96.1%	41.9%
Llama-3.1-8B tuned	Natural Language	87.9%	44.7%
Llama-3.1-8B tuned	JSON	92.3%	46.6%
Llama-3.1-8B tuned	Python Static	87.2%	52.8%
Llama-3.1-8B tuned	Python Dynamic	96.4%	59.2%

Table 8: LLM Two Hop Reasoning with Input Context

cases we studied. In some cases, the fine-tuned small LLMs perform better than much larger LLMs, even though the deep neural

networks of the larger LLMs can help the model make connections

among multiple hops and infer the correct answers. Because we guide and fine-tune the models with emphasis on multi-hop relationships rather than the facts of individual entities, the fine-tuned models can easily generalize to do multi-hop reasoning for completely different entities, even in cases where the number of hops is greater than what is in the training data. Our proposed fine-tuning approaches also help improve the results of in-context learning.

As synthetic data becomes increasingly important for LLM pre-training and fine-tuning, generating proper representations of structured data (especially knowledge graphs) and incorporating this type of data in LLM pre-training and fine-tuning can greatly improve LLM reasoning abilities and reduce hallucinations. As demonstrated by our experimental results, the representation of entity relationships is very important for LLM performance. Programming languages provide the flexibility to represent various entity relationships with native data structures. They can guide LLM reasoning in a more controlled and principled way and ground LLM inference with a knowledge base.

In our experiments, we only studied the simplest form of reasoning: two-hop and three-hop reasoning with compositional relationships. The reasoning process for real-life applications can be much more sophisticated. As shown in Figure 4, the Python representation of a knowledge graph with dynamic relationships provides the flexibility to define any relationships, even extending to a subgraph. Entities and relationships can be defined as classes themselves, allowing us to add attributes to the entities and relationships. Correspondingly, the “infer” function can be redefined to include code that checks these attributes. In future work, we will study the graph representation of entity relationships and its impact on more complex reasoning cases.

5 CONCLUDING REMARKS

We proposed different representations of entity relationships in knowledge graphs to improve the multi-hop reasoning capabilities of LLMs. We conducted a series of experiments to study how different representations of entity relationships affect LLM reasoning ability. We showed that introducing programming language representations of the entity relationships helps improve LLM multi-hop reasoning ability and reduce hallucination.

The programming language representation of the entity relationships provides a controlled and unambiguous way for LLM multi-hop reasoning. By leveraging the native data structures inherent in programming languages, we can effectively model complex entity relationships, while iterative inference functions guide the logical reasoning process. This approach not only enhances reasoning accuracy but also facilitates generalization to more sophisticated reasoning use cases. However, accurately measuring the performance of LLM reasoning beyond two-hop and three-hop compositional relationships can be challenging, as the reasoning process becomes increasingly complex. As part of the future work, we would like to study the programming language representations of more sophisticated relationships in order to solve more complex reasoning tasks. We will experiment at both the pre-training and fine-tuning stages of LLMs to evaluate the performance impact. We hope that our work will inspire other researchers to further push the frontier of LLM research.

REFERENCES

- [1] Meta AI. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] <https://arxiv.org/abs/2407.21783>
- [2] Viraat Aryabumi, Yixuan Su, Raymond Ma, Adrien Morisot, Ivan Zhang, Acyr Locatelli, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. 2024. To Code, or Not To Code? Exploring Impact of Code in Pre-training. arXiv:2408.10914 [cs.CL] <https://arxiv.org/abs/2408.10914>
- [3] Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva, and Amir Globerson. 2024. Hopping Too Late: Exploring the Limitations of Large Language Models on Multi-Hop Queries. arXiv:2406.12775 [cs.CL] <https://arxiv.org/abs/2406.12775>
- [4] Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD Conference*. <https://api.semanticscholar.org/CorpusID:207167677>
- [5] Tim Bray. 2014. The javascript object notation (json) data interchange format.
- [6] Felix Brei, Johannes Frey, and Lars-Peter Meyer. 2024. Leveraging small language models for Text2SPARQL tasks to improve the resilience of AI assistance. arXiv:2405.17076 [cs.AI] <https://arxiv.org/abs/2405.17076>
- [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. arXiv:2005.14165 [cs.CL] <https://arxiv.org/abs/2005.14165>
- [8] Diego Bustamante and Hideaki Takeda. 2024. SPARQL Generation with Entity Pre-trained GPT for KG Question Answering. arXiv:2402.00969 [cs.CL] <https://arxiv.org/abs/2402.00969>
- [9] Ziwei Chai, Tianjie Zhang, Liang Wu, Kaiqiao Han, Xiaohai Hu, Xuanwen Huang, and Yang Yang. 2023. GraphLLM: Boosting Graph Reasoning Ability of Large Language Model. arXiv:2310.05845 [cs.CL] <https://arxiv.org/abs/2310.05845>
- [10] Nurendra Choudhary, Nikhil Rao, Karthik Subbian, and Chandan Reddy. 2022. Graph-based multilingual language model: Leveraging product relations for search relevance. In *KDD 2022*. <https://www.amazon.science/publications/graph-based-multilingual-language-model-leveraging-product-relations-for-search-relevance>
- [11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shrivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways. arXiv:2204.02311 [cs.CL] <https://arxiv.org/abs/2204.02311>
- [12] Stefan Dernbach, Khushbu Agarwal, Alejandro Zuniga, Michael Henry, and Sutanay Choudhury. 2024. GLaM: Fine-Tuning Large Language Models for Domain Knowledge Graph Alignment via Neighborhood Partitioning and Generative Subgraph Encoding. arXiv:2402.06764 [cs.AI] <https://arxiv.org/abs/2402.06764>
- [13] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. arXiv:2404.16130 [cs.CL] <https://arxiv.org/abs/2404.16130>
- [14] Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2023. Complexity-Based Prompting for Multi-step Reasoning. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=yf1icZHC-l9>
- [15] Jiayan Guo, Lun Du, Hengyu Liu, Mengyu Zhou, Xinyi He, and Shi Han. 2023. GPT4Graph: Can Large Language Models Understand Graph Structured Data? An Empirical Evaluation and Benchmarking. arXiv:2305.15066 [cs.AI] <https://arxiv.org/abs/2305.15066>
- [16] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Barcelona, Spain (Online), 6609–6625. <https://www.aclweb.org/anthology/2020.coling-main.580>
- [17] 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 *International Conference on Learning Representations*. <https://openreview.net/forum?id=nZeVKeFYf9>

- [18] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of Hallucination in Natural Language Generation. *Comput. Surveys* 55, 12 (March 2023), 1–38. <https://doi.org/10.1145/3571730>
- [19] Seungone Kim, Se June Joo, Doyoung Kim, Joel Jang, Seonghyeon Ye, Jamin Shin, and Minjoon Seo. 2023. The CoT Collection: Improving Zero-shot and Few-shot Learning of Language Models via Chain-of-Thought Fine-Tuning. arXiv:2305.14045 [cs.CL] <https://arxiv.org/abs/2305.14045>
- [20] Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internet-augmented language models through few-shot prompting for open-domain question answering. arXiv:2203.05115 [cs.CL] <https://arxiv.org/abs/2203.05115>
- [21] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, and Christian Bizer. 2014. DBpedia - A Large-scale, Multilingual Knowledge Base Extracted from Wikipedia. *Semantic Web Journal* 6 (01 2014). <https://doi.org/10.3233/SW-140134>
- [22] Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2024. Reasoning on Graphs: Faithful and Interpretable Large Language Model Reasoning. arXiv:2310.01061 [cs.CL] <https://arxiv.org/abs/2310.01061>
- [23] Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. 2023. At Which Training Stage Does Code Data Help LLMs Reasoning? arXiv:2309.16298 [cs.CL] <https://arxiv.org/abs/2309.16298>
- [24] Zhijie Nie, Richong Zhang, Zhongyuan Wang, and Xudong Liu. 2024. Code-Style In-Context Learning for Knowledge-Based Question Answering. arXiv:2309.04695 [cs.CL] <https://arxiv.org/abs/2309.04695>
- [25] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah Smith, and Mike Lewis. 2023. Measuring and Narrowing the Compositionality Gap in Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 5687–5711. <https://doi.org/10.18653/v1/2023.findings-emnlp.378>
- [26] Julio C. Rangel, Tarcisio Mendes de Farias, Ana Claudia Sima, and Norio Kobayashi. 2024. SPARQL Generation: an analysis on fine-tuning OpenLLaMA for Question Answering over a Life Science Knowledge Graph. arXiv:2402.04627 [cs.AI] <https://arxiv.org/abs/2402.04627>
- [27] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: A Free Collaborative Knowledgebase. *Commun. ACM* 57, 10 (Sept. 2014), 78–85. <https://doi.org/10.1145/2629489>
- [28] Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov. 2024. Can Language Models Solve Graph Problems in Natural Language? arXiv:2305.10037 [cs.CL] <https://arxiv.org/abs/2305.10037>
- [29] Jianing Wang, Junda Wu, Yupeng Hou, Yao Liu, Ming Gao, and Julian McAuley. 2024. InstructGraph: Boosting Large Language Models via Graph-centric Instruction Tuning and Preference Alignment. arXiv:2402.08785 [cs.CL] <https://arxiv.org/abs/2402.08785>
- [30] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=1PLINIMMrw>
- [31] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., 24824–24837. https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf
- [32] Silei Xu, Shicheng Liu, Theo Culhane, Elizaveta Pertseva, Meng-Hsi Wu, Sina Semnani, and Monica Lam. 2023. Fine-tuned LLMs Know More, Hallucinate Less with Few-Shot Sequence-to-Sequence Semantic Parsing over Wikidata. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 5778–5791. <https://doi.org/10.18653/v1/2023.emnlp-main.353>
- [33] Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. 2024. Do Large Language Models Latently Perform Multi-Hop Reasoning? arXiv:2402.16837 [cs.CL] <https://arxiv.org/abs/2402.16837>
- [34] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. In *Thirty-seventh Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=5Xc1ecxO1h>
- [35] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *International Conference on Learning Representations (ICLR)*.
- [36] Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D Manning, Percy Liang, and Jure Leskovec. 2022. Deep Bidirectional Language-Knowledge Graph Pretraining. arXiv:2210.09338 [cs.CL] <https://arxiv.org/abs/2210.09338>
- [37] Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. 2024. Language is All a Graph Needs. arXiv:2308.07134 [cs.CL] <https://arxiv.org/abs/2308.07134>
- [38] Xinlu Zhang, Zhiyu Zoey Chen, Xi Ye, Xianjun Yang, Lichang Chen, William Yang Wang, and Linda Ruth Petzold. 2024. Unveiling the Impact of Coding Data Instruction Fine-Tuning on Large Language Models Reasoning. arXiv:2405.20535 [cs.AI] <https://arxiv.org/abs/2405.20535>
- [39] Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. 2024. Chain of Preference Optimization: Improving Chain-of-Thought Reasoning in LLMs. arXiv:2406.09136 [cs.CL] <https://arxiv.org/abs/2406.09136>
- [40] Kerui Zhu, Bo-Wei Huang, Bowen Jin, Yizhu Jiao, Ming Zhong, Kevin Chang, Shou-De Lin, and Jiawei Han. 2024. Investigating Instruction Tuning Large Language Models on Graphs. arXiv:2408.05457 [cs.CL] <https://arxiv.org/abs/2408.05457>

A ONE SHOT EXAMPLE FOR DIFFERENT PROMPT FORMATS

A.1 Explanation for Natural Language Prompt

The composer of It Goes Like It Goes **is** David Shire. The spouse of David Shire **is** Didi Conn. The spouse of the composer of It Goes Like It Goes **is** _

A.2 JSON structure Prompt

```
{
  "composer": {
    "It_Goes_Like_It_Goes": "David_Shire"
  },
  "spouse": {
    "David_Shire": "Didi_Conn"
  }
}
```

A.3 Python Code Snippet for Prompt V1

```
# Step 1. Define relationships with explicit types
relationships = {
  'composer': {
    'It_Goes_Like_It_Goes': 'David_Shire' # It Goes
    Like It Goes is related to David Shire via
    relationship composer
  },
  'spouse': {
    'David_Shire': 'Didi_Conn' # David Shire is
    related to Didi Conn via relationship spouse
  }
}

# Define entities and relationships
e1 = 'It_Goes_Like_It_Goes'
r1 = 'composer'
r2 = 'spouse'

# Step 2. (r1, e1) -> e2
e2 = relationships[r1][e1]

# Step 3. (r2, e2) -> e3
e3 = relationships[r2][e2]

# Output the result
print(f"{r2}_of_{r1}_is_{e3}")

# when you run the code, it will output:
# The spouse of the composer of It Goes Like It Goes is
Didi Conn
```

A.4 Python Code Snippet for Prompt V2

```
# Step 1. Define relationships with knowledge base
class KnowledgeBase:
    def __init__(self):
        # Initialize an empty dictionary to store facts.
        # Each key is a tuple (entity1, relation),
        # and the value is the entity2 related to
        # entity1 through relation.
        self.facts = {}

    def add_fact(self, entity1, relation, entity2):
        # Add a fact to the knowledge base.
        # :param entity1: The starting entity.
        # :param relation: The relation from entity1
        # to entity2.
        # :param entity2: The related entity reached
        # via the relation.

        self.facts[(entity1, relation)] = entity2

    def infer(self, entity, *relations):
        # Infer the resulting entity by traversing
        # the relations starting from the given
        # entity.

        # :param entity: The starting entity.
        # :param relations: A chain of relations to
        # traverse.
        # :return: The resulting entity after
        # applying the relations, or None if no
        # such path exists.

        current_entity = entity
        for relation in relations:
            key = (current_entity, relation)
            if key in self.facts:
                current_entity = self.facts[key]
            else:
                # If the path does not exist, return
                # None.
                return None
        return current_entity

# Example usage:
# Create a knowledge base instance.
kb = KnowledgeBase()

# Step 2. Define entities and relationships
e1 = 'It_Goes_Like_It_Goes'
r1 = 'composer'
e2 = 'David_Shire'
r2 = 'spouse'
e3 = 'Didi_Conn'

# Add entities and relationships to the knowledge
# base.
kb.add_fact(e1, r1, e2)
kb.add_fact(e2, r2, e3)

# Step 3. Perform inference.
result1 = kb.infer(e1, r1)          # Should return
    David Shire, (It Goes Like It Goes, composer)
    -> David Shire
result2 = kb.infer(e2, r2)          # Should return
    Didi Conn, (David Shire, spouse) -> Didi Conn
```

```
result3 = kb.infer(e1, r1, r2)      # Should return
    Didi Conn, (It Goes Like It Goes, composer) ->
    David Shire, (David Shire, spouse) -> Didi Conn

# Output the result
print(result1) # Output: David Shire is related to
    It Goes Like It Goes through composer
print(result2) # Output: Didi Conn is related to
    David Shire through spouse
print(result3) # Output: Didi Conn is related to It
    Goes Like It Goes through composer and spouse
```