SIMPLE IS EFFECTIVE: THE ROLES OF GRAPHS AND LARGE LANGUAGE MODELS IN KNOWLEDGE-GRAPH-BASED RETRIEVAL-AUGMENTED GENERATION

Mufei Li*, Siqi Miao*, Pan Li Georgia Institute of Technology {mufei.li, siqi.miao, panli}@gatech.edu *Equal contribution.

ABSTRACT

Large Language Models (LLMs) demonstrate strong reasoning abilities but face limitations such as hallucinations and outdated knowledge. Knowledge Graph (KG)-based Retrieval-Augmented Generation (RAG) addresses these issues by grounding LLM outputs in structured external knowledge from KGs. However, current KG-based RAG frameworks still struggle to optimize the trade-off between retrieval effectiveness and efficiency in identifying a suitable amount of relevant graph information for the LLM to digest. We introduce SubgraphRAG, extending the KG-based RAG framework that retrieves subgraphs and leverages LLMs for reasoning and answer prediction. Our approach innovatively integrates a lightweight multilayer perceptron with a parallel triple-scoring mechanism for efficient and flexible subgraph retrieval while encoding directional structural distances to enhance retrieval effectiveness. The size of retrieved subgraphs can be flexibly adjusted to match the query's need and the downstream LLM's capabilities. This design strikes a balance between model complexity and reasoning power, enabling scalable and generalizable retrieval processes. Notably, based on our retrieved subgraphs, smaller LLMs like Llama3.1-8B-Instruct deliver competitive results with explainable reasoning, while larger models like GPT-40 achieve state-of-the-art accuracy compared with previous baselines—all without fine-tuning. Extensive evaluations on the WebQSP and CWQ benchmarks highlight SubgraphRAG's strengths in efficiency, accuracy, and reliability by reducing hallucinations and improving response grounding. Our implementation is available at https://github.com/Graph-COM/SubgraphRAG.

1 Introduction

Large language models (LLMs) have increasingly demonstrated remarkable reasoning capabilities across various domains (Brown et al., 2020; Kojima et al., 2022; Wei et al., 2022; Bubeck et al., 2023; Yao et al., 2023; Huang & Chang, 2023). However, issues like hallucinations (Ji et al., 2023; Huang et al., 2023; Zhang et al., 2023), outdated knowledge (Dhingra et al., 2022; Kasai et al., 2023), and a lack of vertical, domain-specific expertise (Li et al., 2023b) undermine the trustworthiness of LLM outputs. Retrieval-augmented generation (RAG) has emerged as a promising strategy to mitigate these problems by grounding LLM outputs in external knowledge sources (Shuster et al., 2021; Borgeaud et al., 2022; Vu et al., 2024; Gao et al., 2024b).

Despite the effectiveness of text-based retrieval, graph structures offer a more efficient alternative for organizing knowledge (Chein & Mugnier, 2008). Graphs facilitate explicit representation of relationships, reduce information redundancy, and allow for more flexible updates (Robinson et al., 2015). Recent studies have explored using graph-structured knowledge, particularly knowledge graphs (KGs), as external resources for RAG (Pan et al., 2024; Peng et al., 2024; Edge et al., 2024). However, developing effective and efficient frameworks for KG-based RAG remains limited due to the unique challenges involved in retrieving information from the complex structures of KGs.

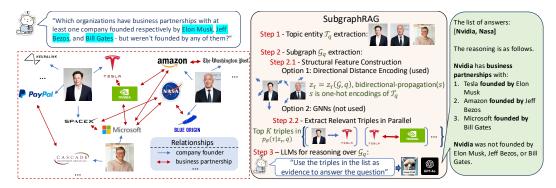


Figure 1: The framework of SubgraphRAG. Retrieved subgraphs consist of relevant triples that are extracted in parallel. Retrieved subgraphs are flexible in their forms and their sizes. In the above example, the relevant subgraph has flexible and complex forms (neither trees nor paths).

Firstly, traditional text-based retrieval methods, such as BM25 (Robertson et al., 1994; Robertson & Zaragoza, 2009) or dense retrieval with cosine similarity (Karpukhin et al., 2020), are insufficient for supporting LLMs in complex reasoning tasks (Sun et al., 2018). For instance, a query like "What is the most famous painting by a contemporary of Michelangelo and Raphael?" requires not only retrieving works by their contemporaries but also reasoning about relationships beyond Michelangelo and Raphael themselves. Thus, KG retrieval goes beyond basic entity linking, requiring the extraction of nonlocal entities connected through multi-hop, relevant relationships to support reasoning (Jiang et al., 2023a; Luo et al., 2024; Sun et al., 2024a). Such information is often best represented as KG subgraphs, whose retrieval enables more effective downstream reasoning.

Second, KG-based RAG faces significant computational challenges. Traditional efficient search methods, such as locality-sensitive hashing, which are designed for similarity search, are not well-suited for extracting complex structural patterns such as paths or subgraphs. With the need to handle potential online graph queries and adapt to dynamic updates in KGs (Trivedi et al., 2017; Liang et al., 2024), efficiently identifying relevant structural information while meeting latency requirements is crucial for designing a practical KG-based RAG framework.

Third, extracting structural information must cover the critical evidence needed to answer the query without exceeding the reasoning capacity of LLMs. Expanding the context window increases computational complexity and can degrade RAG performance by introducing irrelevant information (Xu et al., 2024) and causing the "lost in the middle" phenomenon (Liu et al., 2024b). To prevent these issues, redundant structural information should be pruned to keep only relevant evidence within the LLMs' processing limits, improving accuracy and avoiding hallucinations of LLMs.

Existing KG-based RAG frameworks face limitations in addressing the aforementioned challenges, often due to suboptimal balancing between information retrieval and reasoning over the retrieved data. For example, many approaches rely on LLMs to perform retrieval through step-by-step searches from entities to their neighbors over KGs, resulting in significant complexity by requiring multiple LLM calls (e.g., GPT-4) for each query (Kim et al., Gao et al., 2024a; Wang et al., 2024; Guo et al., 2024; Ma et al., 2024; Sun et al., 2024a; Jiang et al., 2024; Jin et al., 2024). These methods may also miss relevant entities or relationships due to the vast search space over KGs and the limited context windows of LLMs. Conversely, methods that employ lighter models for retrieval, such as LSTMs or GNNs, embed iterative reasoning within the retrieval process itself (Zhang et al., 2022; Liu et al., 2024a; Sun et al., 2019). While more efficient, these methods are constrained by the limited reasoning capacity of lighter models, which can lead to the omission of crucial evidence needed to answer queries. Additionally, some approaches retrieve fixed types of subgraphs for efficiency such as paths (Zhang et al., 2022; Luo et al., 2024), but this restricts the coverage of critical evidence needed for LLM reasoning—a point we will explore more in Sec. 3.1.

Design Principles We argue that there is an inherent tradeoff between model complexity and reasoning capability. To effectively search over KGs, which are expected to grow rapidly, knowledge retrievers should remain lightweight, flexible, generalizable, and equipped with basic reasoning abilities to efficiently filter relevant (even if only roughly relevant) information from vast amounts of irrelevant data, while delegating complex reasoning tasks to LLMs. As LLMs continue to demonstrate increasingly sophisticated reasoning capabilities and are likely to improve further, this division of la-

bor becomes more reasonable. As long as the retrieved information fits within the LLM's reasoning capacity, LLMs—leveraging their superior reasoning power—can then perform more fine-grained analysis and provide accurate answers with appropriate prompting. This approach extends the two-stage Recall & Ranking (rough to fine) framework commonly used in the traditional pipelines of information retrieval and recommendation, with the key advance being that each stage is paired with appropriately tiered reasoning capabilities from AI models to meet the demands of responding to complex queries. Besides, for questions proven too challenging, the concept of iterating the above process can be adopted. However, this consideration is beyond the scope of the current work.

Present Work Our KG-based RAG framework, SubgraphRAG (Fig. 1), follows a pipeline that first retrieves a relevant subgraph and then employs LLMs to reason over it. While this approach mirrors some existing methods (He et al., 2021; Jiang et al., 2023b), SubgraphRAG introduces novel design elements that significantly improve both efficiency and effectiveness by adhering to the aforementioned principles. For efficiency, we employ a lightweight multilayer perceptron (MLP) combined with parallel triple-scoring for subgraph retrieval. To ensure effectiveness, we encode tailored structural distances from the topic entities of a query as structural features. This enables our MLP retriever to outperform more complex models, such as GNNs, LLMs, and heuristic searches, in terms of covering the triples and entities critical for answering the query while maintaining high efficiency. Additionally, the retrieved subgraphs have flexible forms, with adjustable sizes to accommodate the varying capacities of LLMs. SubgraphRAG employs unfine-tuned LLMs, maintaining generalization, adaptability to updated KGs, and compatibility with black-box LLMs.

We evaluate SubgraphRAG on two prominent multi-hop knowledge graph question answering (KGQA) benchmarks—WebQSP and CWQ. Remarkably, without fine-tuning, smaller models like Llama3.1-8B-Instruct can achieve competitive performance. Larger models, such as GPT-40, deliver state-of-the-art (SOTA) results, surpassing previous methods for most cases. Furthermore, SubgraphRAG shows robust multi-hop reasoning capabilities, excelling on more complex, multi-hop questions and demonstrating effective generalization across datasets despite domain shifts. Ablation studies on different retrievers highlight the advantage of our retriever, which consistently outperforms baseline retrievers and is key to our superior KGQA performance. Additionally, our method also exhibits a substantial capability of reducing hallucination by generating knowledge-grounded answers and explanations for its reasoning.

2 Preliminaries

A KG can be represented as a set of triples, denoted by $\mathcal{G} = \{(h,r,t) \mid h,t \in \mathcal{E}, r \in \mathcal{R}\}$, where \mathcal{E} represents the set of entities and \mathcal{R} represents the set of relations. Each triple denoted by $\tau = (h,r,t)$ characterizes a fact that the head entity h and the tail entity t follow a directed relation t. In practice, entities and relations are often associated with a raw text surface form friendly for LLM reasoning.

KG-based RAG aims to enhance LLM responses by incorporating knowledge from a KG as contextual information. Given a query q, LLMs can access relevant knowledge represented by triples in the KG to address the request posed by q. The challenge lies in efficiently searching for relevant knowledge within the often large-scale KG and reasoning to generate an accurate response.

Entity Linking Entity linking is often the first step in KG-based RAG, whose goal is to identify the set of entities $\mathcal{T}_q \subset \mathcal{E}$ directly involved in the query q. The entities in \mathcal{T}_q , named topic entities, provide valuable inductive bias for retrieval as the triples relevant to q are often close to \mathcal{T}_q .

Knowledge Graph Question Answering (KGQA) is a key application often used to evaluate KG-based RAG, where the query q is a question that requires finding answers under specific constraints. The answer(s) \mathcal{A}_q typically corresponds to a set of entities in the KG. Questions that require multiple triples in the KG as evidence to identify an answer entity are classified as complex questions, as they demand multi-hop reasoning, in contrast to single-hop questions.

3 THE SUBGRAPHRAG FRAMEWORK

Our proposed framework, SubgraphRAG, adopts a retrieval-and-reasoning pipeline for KG-based RAG. Specifically, given a query q, SubgraphRAG first extracts a subgraph $\mathcal{G}_q \subset \mathcal{G}$ that represents relevant knowledge, and then LLMs generate the response by reasoning over \mathcal{G}_q . While some exist-

ing works follow a similar pipeline (Zhang et al., 2022; Luo et al., 2024; Wu et al., 2023; Wen et al., 2023), our approach is novel in addressing three key challenges: First, the extracted subgraph \mathcal{G}_q is designed to cover as much of the relevant evidence for answering q as possible, while adhering to a size constraint K, which can be adjusted according to the capacity of the downstream LLM. Second, the extraction of \mathcal{G}_q is highly efficient and scalable. Third, we employ tailored prompting to guide the LLM in reasoning over \mathcal{G}_q and generating a well-grounded answer with explanations.

3.1 EFFICIENT, FLEXIBLE AND EXPRESSIVE SUBGRAPH RETRIEVAL

Problem Reduction To begin with, we formulate the subgraph retrieval problem and gradually reduce it to an efficiently solvable problem. An LLM can be viewed as an answer generator $\mathbb{P}(\cdot | \mathcal{G}_q, q)$ that takes queries and evidence represented by subgraphs. Given a query q and its answer \mathcal{A}_q , the best subgraph evidence for this LLM is denoted as $\mathcal{G}_q^* = \arg\max_{\mathcal{G}_q \subseteq \mathcal{G}} \mathbb{P}(\mathcal{A}_q \mid \mathcal{G}_q, q)$. Of course, solving this problem is practically impossible as it requires the knowledge of \mathcal{A}_q . Instead, we aim to learn a subgraph retriever from data and expect this retriever to generalize to unseen future queries.

Specifically, let the subgraph retriever be a distribution $\mathbb{Q}_{\theta}(\cdot|q,\mathcal{G})$ over the subgraph space of the KG. θ denotes the parameters. Given a training set of question-answer pairs \mathcal{D} , the subgraph retriever learning problem can be formulated as the following problem:

$$\max_{\theta} \mathbb{E}_{(q, \mathcal{A}_q) \sim \mathcal{D}, \mathcal{G}_q \sim \mathbb{Q}_{\theta}(\mathcal{G}_q | q, \mathcal{G})} \mathbb{P}(\mathcal{A}_q \mid \mathcal{G}_q, q). \tag{1}$$

In practice, this problem is still hard to solve due to the complexity of the LLM, i.e., the form of \mathbb{P} . Simply evaluating $\mathbb{P}(\mathcal{A}_q \mid \mathcal{G}_q, q)$ means calling the LLM to generate the particular answer \mathcal{A}_q , which could be costly and only applicable to grey/white-box LLMs with accessible output logits, let alone the incomputable gradient $\frac{d\mathbb{P}}{d\mathcal{G}_q}$.

To solve the problem in Eq. 1, we adopt the following idea. If we know the optimal subgraph \mathcal{G}_q^* , the maximum likelihood estimation (MLE) principle can be leveraged to train the retriever $\max_{\theta} \mathbb{E}_{(q,\mathcal{A}_q)\sim\mathcal{D}} \mathbb{Q}_{\theta}(\mathcal{G}_q^* \mid \mathcal{G},q)$. However, getting \mathcal{G}_q^* for an even known question-answer pair (q,\mathcal{A}_q) is computationally hard and LLM-dependent. Instead, we use (q,\mathcal{A}_q) to construct surrogate subgraph evidence with heuristics $\tilde{\mathcal{G}}(q,\mathcal{A}_q)$ and train the retriever based on MLE:

$$\max_{q} \ \mathbb{E}_{(q,\mathcal{A}_q)\sim\mathcal{D}} \ \mathbb{Q}_{\theta}(\tilde{\mathcal{G}}_q \mid \mathcal{G}, q), \quad \text{where } \tilde{\mathcal{G}}_q = \tilde{\mathcal{G}}(q, \mathcal{A}_q).$$
 (2)

Some examples of $\tilde{\mathcal{G}}_q$ could be the shortest paths between topic entities \mathcal{T}_q and the answer entities \mathcal{A}_q . Eq. 2 is conceptually similar to the weak supervision adopted in some existing work (Zhang et al., 2022). However, the formulation in Eq. 2 indicates that the sampled subgraph does not necessarily follow a fixed type (trees or paths). Instead, the retriever distribution \mathbb{Q}_θ can by construction factorize into a product of distributions over triples, allowing efficient training and inference, flexible subgraph forms, and adjustable subgraph sizes.

Triple Factorization We propose to adopt a retriever that allows a subgraph distribution factorization over triples given some latent variables $z_{\tau} = z_{\tau}(\mathcal{G}, q)$ (to be elaborated later):

$$\mathbb{Q}_{\theta}(\mathcal{G}_q \mid \mathcal{G}, q) = \prod_{\tau \in \mathcal{G}_q} p_{\theta}(\tau \mid z_{\tau}, q) \prod_{\tau \in \mathcal{G} \setminus \mathcal{G}_q} (1 - p_{\theta}(\tau \mid z_{\tau}, q)),$$

This strategy is inspired by the studies on graph generative models (Kipf & Welling, 2016) and enjoys four benefits: Efficiency in Training - The problem in Eq. 2 can be factorized as $\max_{\theta} \mathbb{E}_{(q,\mathcal{A}_q)\sim\mathcal{D}} \sum_{\tau\in\tilde{\mathcal{G}}_q} \log p_{\theta}(\tau\mid z_{\tau},q) + \sum_{\tau\in\mathcal{G}\setminus\tilde{\mathcal{G}}_q} \log(1-p_{\theta}(\tau\mid z_{\tau},q));$ Efficiency in Sampling - After computing z_{τ} , we can select triples τ from \mathcal{G} in parallel; Flexibility - Triple combinations can form arbitrary subgraphs; Adjustable Size - Subgraphs formed by top-K triples with different K values can accommodate various LLMs with diverse reasoning capabilities. In practice, \mathbb{Q}_{θ} can be further simplified given topic entities \mathcal{T}_q (He et al., 2021; Jiang et al., 2023b), by only considering subgraphs close to the topic entities, i.e., $p_{\theta}(\tau\mid z_{\tau},q)=0$ for a τ that is far from \mathcal{T}_q .

Relevant Designs Previous approaches often adopt heuristics and focus on some particular types of subgraphs, such as constrained subgraph search (e.g., searching for connected subgraphs (He et al., 2024)), constrained path search from topic entities, often imposing constraints in path counts and lengths and employing expensive iterative processes (Zhang et al., 2022; Wu et al., 2023; Luo et al.,

```
System: Based on the triples retrieved from a knowledge graph, please answer the question. Please return formatted answers as a list, each prefixed with "ans:". User: Triplets: (e_1, r_{12}, e_2) \setminus n (e_3, r_{34}, e_4) \setminus n \dots \setminus n Question: ... // ICL example Assistant: To answer the question, we have to find .... From the triples we can see that .... Therefore, the answers are: \setminus n ans: ... \setminus n ans: ... \setminus n .... // ICL example User: Triplets: (e_a, r_{ab}, e_b) \setminus n (e_c, r_{cd}, e_d) \setminus n ... \setminus n Question: ... // the evaluation question Assistant: To answer the question, we have to find .... From the triples we can see that .... Therefore, the answers are: \setminus n ans: ... \setminus n answer
```

Figure 2: The prompt used in SubgraphRAG. Concrete examples can be found in Appendix D).

2024; Sun et al., 2024a; Liu et al., 2024a; Mavromatis & Karypis, 2024; Sun et al., 2024b), and entity selection followed by extracting entity-induced subgraphs, where all triples involving an entity are included together (Yasunaga et al., 2021; Taunk et al., 2023). The loss in flexibility narrows the space of possible retrieved subgraphs, which eventually harms the effectiveness of the RAG.

Directional Distance Encoding (DDE) as $z_{\tau}(\mathcal{G},q)$ The latent variable $z_{\tau}(\mathcal{G},q)$ aims to model the relationship between a triple τ and the query q given \mathcal{G} . One idea is to employ graph neural networks (GNNs) to compute $z_{\tau}(\mathcal{G},q)$ through message passing between entities/relations with attribute embeddings and question embeddings. Some previous works indeed adopt GNNs to get latent representations of entities (Yasunaga et al., 2021; Kang et al., 2023; Mavromatis & Karypis, 2024; Liu et al., 2024a). However, GNNs are known to have limited representation power (Xu et al., 2019; Morris et al., 2019; Chen et al., 2020).

The structural relationship between τ and q provides valuable information complementing to their semantic relationship. Inspired by the success of distance encoding and labeling trick in enhancing the structural representation power of GNNs (Li et al., 2020; Zhang et al., 2021), we propose a DDE as $z_{\tau}(\mathcal{G},q)$ to model the structural relationship. Given topic entities \mathcal{T}_q , let $\mathbf{s}_e^{(0)}$ be a one-hot encoding representing $e \in \mathcal{T}_q$ or $e \notin \mathcal{T}_q$. For the l+1-th round, we perform feature propagation and compute $\mathbf{s}_e^{(l+1)} = \text{MEAN}\{\mathbf{s}_{e'}^{(l)} \mid (e',\cdot,e) \in \mathcal{G}\}$, and through the reverse direction to account the directed nature of \mathcal{G} , $\mathbf{s}_e^{(r,l+1)} = \text{MEAN}\{\mathbf{s}_{e'}^{(r,l)} \mid (e,\cdot,e') \in \mathcal{G}\}$, where $\mathbf{s}_e^{(r,0)} = \mathbf{s}_e^{(0)}$. We concatenate the results across all rounds and both directions to obtain the final entity encodings $\mathbf{s}_e = [\mathbf{s}_e^{(0)} \| \mathbf{s}_e^{(1)} \| \cdots \| \mathbf{s}_e^{(r,1)} \| \cdots]$, which leads to triple encodings as $z_{\tau}(\mathcal{G},q) = [\mathbf{s}_h \| \mathbf{s}_t]$ that concatenates the head h's and the tail t's encodings. In section 4.1, we compare different approaches to compute $z_{\tau}(\mathcal{G},q)$ - using GNNs, DDEs or only one-hot encodings of \mathcal{T}_q , and DDEs perform the best.

A Lightweight Implementation For $p_{\theta}(\cdot|z_{\tau}(\mathcal{G},q),q)$ We present a lightweight implementation of p_{θ} that integrates structural and semantic information. Following previous approaches (Karpukhin et al., 2020; Gao et al., 2024b), we employ off-the-shelf pre-trained text encoders to embed all entities/relations in a KG based on their text attributes. These semantic text embeddings are computed and stored in a vector database during the pre-processing stage for efficient retrieval. For a newly arrived question q, we embed q to obtain z_q and retrieve embeddings z_h, z_r, z_t from the vector database for the involved entities and relation. After computing DDEs z_{τ} , an MLP is employed for binary classification using the concatenated input $|z_q||z_h||z_r||z_t||z_\tau|$.

Relevant Designs We considered several alternative design options but found them less suitable due to concerns regarding efficiency and adaptability to KG updates. Cross-encoders, which concatenate a question and a retrieval candidate for joint embedding (Wolf et al., 2019), potentially offer better retrieval performance. However, due to the inability to pre-compute embeddings, this approach significantly reduces retrieval efficiency when dealing with a large number of retrieval candidates, as is the case in triple retrieval. Li et al. (2023a) embeds each triple as a whole rather than individual entities and relations. However, this approach incurs higher computational and storage costs and exhibits reduced generalizability to the triples that are new combinations of old entities and relations. Our implementation allows for fast triple scoring while maintaining good generalizability.

3.2 PROMPTING-BASED LLM REASONING

We utilize an LLM to reason over \mathcal{G}_q by incorporating a linearized list of triples from \mathcal{G}_q into the prompt. This enables the LLM to ground its reasoning in the retrieved subgraph and identify the answers $\hat{\mathcal{A}}_q$ from the entities within \mathcal{G}_q , addressing issues such as hallucinations and outdated knowledge (Lin et al., 2019; Shuster et al., 2021; Vu et al., 2024). Specifically, we prompt the LLM not only to provide answers but also to generate knowledge-grounded explanations based on the input

Table 1: Evaluation results for retrieval recall and wall-clock time. Best results are in **bold**. Being training-free, cosine similarity and G-Retriever stay unchanged in generalization evaluations.

Model		WebQ	SP			CWQ)		CWQ	→WebQSI	P	$WebQSP \rightarrow CWQ$		
	Triple	s	Entites		Triple	Triples Entite			Triple	S	Entites	Triple	s	Entites
	Shortest Path	GPT-40	Answer	Time (s)	Shortest Path	GPT-40	Answer	Time (s)	Shortest Path	GPT-40	Answer	Shortest Path	GPT-40	Answer
cosine similarity	0.714	0.719	0.708	3	0.488	0.567	0.582	13	0.714	0.719	0.708	0.488	0.567	0.582
Retrieve-Rewrite-Answer	0.058	0.062	0.740	69	-	-	-	-	-	-	-	-	-	-
RoG	0.713	0.388	0.807	948	0.623	0.298	0.841	2327	0.589	0.323	0.658	0.301	0.139	0.412
G-Retriever	0.294	0.325	0.545	672	0.183	0.217	0.375	1530	0.294	0.325	0.545	0.183	0.217	0.375
SubgraphRAG	0.883	0.865	0.944	6	0.811	0.840	0.914	<u>12</u>	0.794	0.776	0.887	0.622	0.623	0.773

Table 2: Breakdown of recall evaluation over # hops. Best results are in **bold**.

Model		Т	riple Recall				GPT-	40 Triple R	ecall			Answ	er Entity R	ecall	
	Web	QSP		CWQ		Web	QSP		CWQ		Web	QSP		CWQ	
	1 (65.8%)	2 (34.2%)	1 (28.0%)	2 (65.9%)	≥ 3 (6.1%)	1 (65.8%)	2 (34.2%)	1 (28.0%)	2 (65.9%)	≥ 3 (6.1%)	1 (65.8%)	2 (34.2%)	1 (28.0%)	2 (65.9%)	≥ 3 (6.1%)
cosine similarity	0.874	0.405	0.629	0.442	0.333	0.847	0.483	0.629	0.511	0.464	0.943	0.253	0.903	0.472	0.289
Retrieve-Rewrite-Answer	0.064	0.046	-	-	-	0.062	0.061	-	-	-	0.745	0.729	-	-	-
RoG	0.869	0.415	0.766	0.597	0.253	0.446	0.271	0.347	0.293	0.122	0.874	0.677	0.920	0.827	0.628
G-Retriever	0.335	0.216	0.134	0.205	0.168	0.345	0.284	0.159	0.240	0.226	0.596	0.446	0.377	0.384	0.269
MLP	0.828	0.687	0.651	0.690	0.534	0.811	0.781	0.635	0.707	0.616	0.933	0.874	0.932	0.870	0.793
MLP + topic entity SubgraphRAG	0.944 <u>0.953</u>	0.729 <u>0.748</u>	$\frac{0.854}{0.831}$	0.750 <u>0.820</u>	0.560 <u>0.626</u>	0.884 <u>0.908</u>	0.775 <u>0.809</u>	0.769 <u>0.823</u>	0.773 <u>0.860</u>	0.647 0.755	0.976 <u>0.977</u>	0.843 <u>0.881</u>	0.956 0.946	0.885 <u>0.916</u>	$0.665 \\ 0.741$

subgraph. We adopt in-context learning (ICL) (Brown et al., 2020) and design dedicated prompt templates with explanation demonstrations to guide the LLM's reasoning process (see Fig. 2).

By avoiding the need for fine-tuning LLMs, we reduce computational costs, enable the use of SOTA black-box LLMs, and maintain the framework's generalizability, even for unseen KGs. While fine-tuning may enhance prediction accuracy, it often diminishes general reasoning and explanatory capabilities. Furthermore, high-quality labels for text-based explanations are typically unavailable in practical question-answering tasks. Consequently, previous KG-based RAG approaches that rely on fine-tuning often generate reasoning explanations using larger, unfine-tuned LLMs, such as GPT-4, to serve as auxiliary labels for additional training (Luo et al., 2024).

Regarding the size K of the retrieved subgraph, while increasing K in principle improves the coverage of relevant information, it also incurs higher costs/latency for LLM reasoning and risks introducing more irrelevant information that may ultimately hurt LLM reasoning (Xu et al., 2024; Liu et al., 2024b). Different LLMs are inherently equipped with different sized context window and also exhibit distinct capabilities in reasoning over long-context retrieval results (Dubey et al., 2024; Leng et al., 2024). As such, although the training of SubgraphRAG retriever is LLM-agnostic, the size K needs to be properly selected per LLM and cost/latency constraint. In Section 4.2, we empirically verify that more powerful LLMs can benefit from incorporating a larger-sized retrieved subgraph, demonstrating the benefit of size-adjustable subgraph retrieval in SubgraphRAG.

3.3 MORE RELATED WORKS

Along with the introduction of components in SubgraphRAG, we have introduced the most relevant works. Other related works are discussed in Appendix A due to the space limitation.

4 EXPERIMENTS

We design our empirical studies to examine the effectiveness and efficiency of SubgraphRAG in addressing the various challenges inherent to KG-based RAG, covering both retrieval and reasoning aspects. Q1) Overall, to meet the accuracy and low-latency requirements of KG-based RAG, does SubgraphRAG effectively and efficiently retrieve relevant information? Q2) For complex questions involving multi-hop reasoning and multiple topic entities, does SubgraphRAG properly integrate structural information for effective retrieval? Q3) How effectively does SubgraphRAG perform on KGQA tasks, and how is its accuracy influenced by different factors? Q4) To what extent can our pipeline provide effective knowledge-grounded explanations for question answering?

Datasets. We adopt two prominent and challenging KGQA benchmarks that necessitate multi-hop reasoning – WebQSP (Yih et al., 2016) and CWQ (Talmor & Berant, 2018). Both benchmarks utilize Freebase (Bollacker et al., 2008) as the underlying KG. To evaluate the capability of LLM reasoners in knowledge-grounded hallucination-free question answering, we introduce WebQSP-sub and CWQ-sub, where we remove samples whose answer entities are absent from the KG.

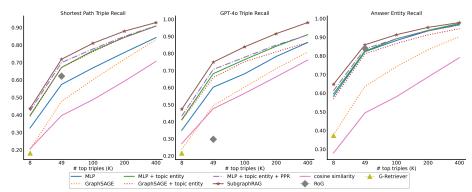


Figure 3: Retrieval effectiveness on CWQ across a spectrum of K values for top-K triple retrieval.

4.1 EVALUATION FOR RETRIEVAL (Q1 & Q2)

Baseline Retrievers. Li et al. (2023a) introduces a structure-free retriever that performs cosine similarity search based on triple embeddings, which we refer to as **cosine similarity**. **Retrieve-Rewrite-Answer** (Wu et al., 2023) proposes a constrained path search, predicting relation paths then searching for matched paths. **RoG** (Luo et al., 2024) adopts a similar strategy but enhances it by fine-tuning an LLM for generative relation path prediction. **G-Retriever** (He et al., 2024) combines cosine similarity search with combinatorial optimization to construct a connected subgraph.

Implementation Details. We employ gte-large-en-v1.5 (Li et al., 2023c) as the pre-trained text encoder for both the cosine similarity baseline and SubgraphRAG, a 434M model that achieves a good balance between efficiency and English retrieval performance, as evidenced by the Massive Text Embedding Benchmark (MTEB) leaderboard (Muennighoff et al., 2023). For supervision signals, there are no ground-truth relevant subgraphs for a query q. Previous path-based subgraph retrievers adopt the shortest paths between the topic and answer entities for the weak supervision signals (Zhang et al., 2022; Luo et al., 2024). We also utilize these shortest paths as the heuristic relevant subgraphs $\tilde{\mathcal{G}}_q$ to train \mathbb{Q}_θ as in Eq. 2. To reduce the size of candidate triples \mathcal{G} , we construct subgraphs centered at topic entities \mathcal{T}_q , following previous works. See Appendix B.1 for more details.

Evaluation Metrics. Our retrieval evaluation encompasses both effectiveness and efficiency. For effectiveness, we employ three recall metrics: recall of triples in shortest paths (i.e., $\tilde{\mathcal{G}}_a$), recall of GPT-40identified relevant triples, and recall of answer entities within retrieved subgraphs or triples. The first metric assesses the ability of the approaches to retrieve the heuristic signals used as weak supervision during training. To provide a more accurate assessment for relevant triple retrieval, we employ GPT-40 to identify up to 20 high-quality relevant triples for all test set samples, potentially capturing relevant triples beyond the shortest paths (see Appendix D for more details on labeling). We calculate individual recall values per question and average across all questions. For efficiency, we measure wall clock time on a 48GB NVIDIA RTX 6000 Ada GPU, reporting results in seconds. To focus on model computational efficiency, we exclude KG query times from our measurements. For text embedding computation, we only consider the time taken for ques-

Table 3: Question-answering performance on WebQSP and CWQ. Best results are in **bold**. By default, our reasoners use the top 100 retrieved triples. Results with 200 and 500 triples (indicated in parentheses) are also shown. Results with (\leftrightarrow) evaluate retriever generalizability, where the retriever is trained on one dataset and applied to the other.

WebQ	SP	CWO	5
Macro-F1	Hit	Macro-F1	Hit
52.5	68.6	-	55.7
-	82.6	-	67.6 ¹
-	82.5	-	62.0
-	74.69	-	-
-	79.36	-	-
53.41	73.46	-	-
70.26	86.67	54.63	61.94
66.45	82.19	53.87	60.55
70.57	86.61	47.16	56.98
74.70	86.24	51.78	57.89
69.21	83.11	49.13	56.27
77.45	90.11	54.13	62.02
76.46	89.80	59.08	66.69
77.82	90.54	54.69	63.49
78.24	90.91	59.42	67.49
77.67	91.22	55.41	64.97
66.42	83.42	37.96	48.57
73.81	88.08	44.69	54.21
76.20	91.22	50.30	60.80
	Macro-F1 52.5 53.41 70.26 66.45 70.57 74.70 69.21 77.45 76.46 77.82 78.24 77.67 66.42 73.81	52.5 68.6 - 82.6 - 82.5 - 74.69 - 79.36 53.41 73.46 70.26 86.67 66.45 82.19 70.57 86.61 74.70 86.24 69.21 83.11 76.46 89.80 77.82 90.54 77.82 90.54 77.67 91.22 66.42 83.42 73.81 88.08	Macro-F1 Hit Macro-F1 52.5 68.6 - - 82.6 - - 82.5 - - 74.69 - - 79.36 - 53.41 73.46 - 70.26 86.67 54.63 66.45 82.19 53.87 70.57 86.61 47.16 74.70 86.24 51.78 69.21 83.11 49.13 77.45 90.11 54.13 76.46 89.80 59.08 77.82 90.54 54.69 77.67 91.22 55.41 66.42 83.42 37.96 73.81 88.08 44.69

tion embedding computation as the entity and relation embeddings for the entire KG can be precomputed.

Overall Evaluation (Q1). Table 1 shows that SubgraphRAG consistently outperforms all other approaches in retrieval effectiveness. RoG achieves the most competitive baseline performance for retrieving shortest path triples and answer entities. However, when evaluated on GPT-40-labeled

triples, RoG's performance drops significantly (45.6% for WebQSP, 52.2% for CWQ), while SubgraphRAG remains robust (2.0% decrease for WebQSP, 3.6% increase for CWQ). This stark difference, despite using the same training signals, empirically validates that SubgraphRAG's individual triple selection mechanism allows for more flexible and effective subgraph extraction compared to RoG's constrained path search approach. Regarding **efficiency**, SubgraphRAG is only slightly slower than the cosine similarity baseline on WebQSP while being one to two orders of magnitude faster than other baselines. For the cosine similarity baseline and SubgraphRAG, we report recall metrics based on the top-100 retrieved triples (2.3% of total candidate triples on average). This budget consistently yields robust reasoning performance across LLMs in subsequent experiments.

Generalizability. We further examine the generalizability of the retrievers by training them on dataset A and evaluating on dataset B, denoted as $A \to B$ in Table 1. Despite an anticipated performance degradation, SubgraphRAG consistently outperforms the alternative approaches.

Ablation Study for Design Options and Retrieval Size. To evaluate individual component contributions in SubgraphRAG, we conduct an ablation study with several variants. MLP is a structure-free variant employing only text embeddings. Given the prevalence of GNNs, we consider Graph-SAGE (Hamilton et al., 2017), a popular GNN, to update entity representations prior to the MLP-based triple scoring. We further augment both MLP and GraphSAGE with a one-hot-encoding topic entity indicator (MLP + topic entity and GraphSAGE + topic entity). Following Sun et al. (2018), we incorporate Personalized PageRank (PPR) (Haveliwala, 2002), seeded from the topic entities, to integrate structural information (MLP + topic entity + PPR). To account for both the varying capabilities of downstream LLMs and inference cost/latency constraints, we evaluate these variants across a broad spectrum of retrieval sizes (K).

Fig. 3 presents the results on CWQ, with baselines included for reference. Larger retrieval sizes uniformly improve recall across all variants. Equipped with DDE, SubgraphRAG outperforms other variants, even at the relatively small average retrieval sizes of the baselines. This demonstrates that SubgraphRAG's superiority is not solely attributable to larger retrieval sizes. Regarding design options, the topic entity indicator invariably leads to an improvement. In contrast, GNN variants often result in performance degradation compared to their MLP counterparts. We suspect that the diffusion of semantic information introduces noise in triple selection. Finally, PPR fails to reliably yield improvements. For the results on WebQSP, see Appendix B.3, which are also consistent.

Multi-Hop and Multi-Topic Questions (Q2). To evaluate the effectiveness of various approaches in capturing structural information for complex multi-hop and multi-topic questions, we group questions based on the number of hops and topic entities. Table 2 presents the performance breakdown by hop count. SubgraphRAG consistently outperforms other methods on WebQSP and achieves the best overall performance on CWQ. Notably, while the cosine similarity baseline and RoG demonstrate competitive performance for single-hop questions, their performance degrades significantly for multi-hop questions. Appendix B.4 provides a performance breakdown for single-topic and multi-topic questions, focusing exclusively on CWQ due to the predominance of single-topic questions in the WebQSP test set (98.3%). SubgraphRAG consistently exhibits superior performance across all metrics for both single-topic and multi-topic questions. Our comprehensive analysis high-lights the remarkable effectiveness of DDE in capturing complex topic-centered structural information essential for challenging questions involving multi-hop reasoning and multiple topic entities.

4.2 KGQA RESULTS (**Q3 & Q4**)

KGQA Baselines. Besides the baselines used for retriever evaluation, we include results from other LLM-based KGQA methods due to their state-of-the-art performance, such as KD-CoT (Wang et al., 2023a), StructGPT (Jiang et al., 2023a), ToG (Sun et al., 2024a), and EtD (Liu et al., 2024a). For RoG, we present two entries: RoG-Joint and RoG-Sep. Originally, RoG fine-tuned its LLMs on the training sets of both WebQSP and CWQ. Yet, this joint training approach leads to significant label leakage, with over 50% of WebQSP test questions (or their variants) appearing in CWQ's training set, and vice versa. Therefore, we re-trained RoG on each dataset separately, indicated as RoG-Sep.

Evaluation Metrics. Along with the commonly reported Macro-F1 and Hit², we also include Micro-F1 to account for the imbalance in the number of ground-truth answers across samples and Hit@1

¹Their computation of Hit is different from other baselines and may overestimate the performance. We were unable to reproduce their results following their provided instructions.

Table 4: Question-answering performance on WebQSP-sub and CWQ-sub. Best results are in **bold**. By default, our reasoners use the top 100 retrieved triples. Results with 200 and 500 triples (indicated in parentheses) are also shown. Results with (\leftrightarrow) evaluate retriever generalizability, where the retriever is trained on one dataset and applied to the other.

		Web	QSP-sub				CW	/Q-sub		
	Macro-F1	Micro-F1	Hit	Hit@1	$Score_h$	Macro-F1	Micro-F1	Hit	Hit@1	Score
G-Retriever	54.13	23.84	74.52	67.56	67.97	-	-	-	-	-
RoG-Joint	72.01	47.70	88.90	82.62	76.13	58.61	52.12	66.22	61.17	55.15
RoG-Sep	67.94	43.10	84.03	77.61	72.79	57.69	52.83	64.64	60.64	54.51
SubgraphRAG + Llama3.1-8B	72.10	46.56	88.58	84.80	82.42	54.76	51.76	65.80	59.69	62.89
SubgraphRAG + Llama3.1-70B	75.97	51.64	87.88	85.89	85.57	61.49	59.91	68.43	65.52	67.62
SubgraphRAG + ChatGPT	70.81	44.73	85.18	80.82	81.53	56.37	54.44	64.40	60.99	61.31
SubgraphRAG + GPT-4o-mini	78.34	58.44	91.34	87.36	82.21	61.13	58.86	70.01	65.48	64.20
SubgraphRAG + GPT-4o	77.61	56.78	91.40	86.40	81.85	65.99	<u>63.18</u>	73.91	68.89	66.57
SubgraphRAG + GPT-4o-mini (200)	78.66	58.65	91.73	87.04	81.98	61.58	57.47	71.45	65.87	63.66
SubgraphRAG + GPT-4o (200)	79.40	58.91	92.43	87.75	82.46	66.48	61.30	75.14	69.42	66.45
SubgraphRAG + GPT-4o-mini (500)	78.46	57.08	92.43	<u>88.01</u>	81.95	62.18	56.86	72.82	66.57	62.77
SubgraphRAG + Llama3.1-8B (↔)	67.91	42.79	85.25	81.21	80.09	43.03	40.73	55.09	47.58	56.78
SubgraphRAG + GPT-4o-mini (↔)	74.42	49.41	89.10	84.67	81.35	49.47	45.16	60.18	54.18	58.86
SubgraphRAG + GPT-4o-mini (↔, 500)	76.83	52.01	92.30	87.43	81.41	55.58	49.13	67.24	59.69	59.87

Table 5: Breakdown of QA performance by reasoning hops.

		WebQ	SP-sub				CWQ-	sub		
	1 (65.8°	(65.8%)		(34.2%)		(28.0%)		%)	≥ 3 (6.1%	
	Marco-F1	Hit	Marco-F1	Hit	Marco-F1	Hit	Marco-F1	Hit	Marco-F1	Hit
G-Retriever	56.41	78.20	45.73	65.35	-	-	-	-	-	-
RoG-Joint	77.05	92.96	62.53	81.54	59.75	66.33	59.70	68.56	41.46	43.27
RoG-Sep	74.50	89.83	55.62	73.45	59.35	66.20	59.45	67.17	31.33	33.33
SubgraphRAG (Llama3.1-8B)	75.50	91.40	65.87	83.62	51.54	63.05	57.52	68.93	41.88	47.37
SubgraphRAG (GPT-4o-mini)	80.56	92.86	74.11	88.51	57.36	67.34	63.85	72.74	51.14	54.39

for a more inclusive evaluation. To further assess model performance, we introduce $score_h$, inspired by (Yang et al., 2024), which evaluates how truth-grounded the predicted answers are and the degree of hallucination. This metric penalizes hallucinated answers while favoring missing answers over incorrect ones. Scores are normalized to a range of 0 to 100, with higher scores indicating better truth-grounding in the model's answers. Details of the scoring strategy are provided in Appendix C.

Experiment Settings. We use the vllm (Kwon et al., 2023) framework for efficient LLM inference. For KD-CoT, Retrieve-Rewrite-Answer, ToG, and EtD, we report their published results due to the difficulty in reproducing them. For StructGPT, we directly use their provided processed files to obtain results for WebQSP. For the remaining baselines, we successfully reproduced their results and evaluated them on additional metrics and datasets, including WebQSP-sub and CWQ-sub. However, we do not include results for G-Retriever on CWQ, as it required over 200 hours of computation on 2 NVIDIA RTX 6000 Ada GPUs. If not specified, the LLM reasoners in SubgraphRAG use the top 100 retrieved triples; results using more triples are explicitly noted. All Llama variants considered are based on instruction tuning. For the LLM reasoners used in SubgraphRAG, both the temperature and the seed are set to 0 to ensure reproducibility.

Overall Performance. Tables 3 and 4 present the evaluation results, where SubgraphRAG achieves state-of-the-art (SOTA) results on both WebQSP and WebQSP-sub. Even with smaller 8B LLMs, our method surpasses previous SOTA approaches by up to 4% in Macro-F1 and Hit metrics (excluding RoG-Joint due to test label leakage in that model). With larger models like Llama3.1-70B-Instruct and GPT-40, SubgraphRAG achieves even greater performance, showing up to a 12% improvement in Macro-F1 and a 9% increase in Hit. On the more challenging CWQ and CWQ-sub datasets, which require extended reasoning hops, SubgraphRAG performs competitively even with smaller 8B models. When paired with advanced reasoning models like GPT-40, SubgraphRAG achieves results second only to ToG on CWQ. Notably, SubgraphRAG requires only a single call to GPT-40, whereas ToG requires 6-8 calls, which increases computational cost, and we were unable to reproduce ToG's performance using their published code. On CWQ-sub, SubgraphRAG demonstrates gains of up to 9% in Macro-F1 and 11% in Hit, indicating that tasks with greater reasoning complexity benefit

²The baselines claim to report Hit@1, but they actually compute Hit, which measures whether at least one correct answer appears in the LLM response.

Table 6: Detailed performance of truth-grounded QA. No Ans Samples refers to cases where LLM reasoners refuse to answer; NR (Not Retrieved) indicates answers not present in the retrieved triples, while R (Retrieved) indicates answers found within the retrieved triples.

Dataset	Method		Samples w/ An	s Entities in KG		Samı	oles w/o Ans Entities	in KG
Dutuset	Method	No Ans Samples Total Samples	Wrong Ans Total Ans	Correct Ans (NR) Correct Ans	Wrong Ans (NR) Wrong Ans	No Ans Samples Total Samples	Wrong Ans (NR) Total Ans	Wrong Ans (R) Total Ans
WebQSP	RoG-Joint	0/1559 = 0%	4605/10206 = 45%	170/5601 = 3%	1178/4605 = 26%	0/69 = 0%	24/126 = 19%	70/126 = 56%
	RoG-Sep	0/1559 = 0%	6867/12401 = 55%	327/5534 = 6%	2494/6867 = 36%	0/69 = 0%	71/162 = 44%	43/162 = 27%
weeger	SubgraphRAG (Llama3.1-8B)	29/1559 = 2%	1397/5850 = 24%	59/4453 = 1%	84/1397 = 6%	13/69 = 19%	16/107 = 15%	71/107 = 66%
	SubgraphRAG (GPT4o-mini)	12/1559 = 1%	1802/8011 = 22%	37/6209 = 1%	126/1802 = 7%	7/69 = 10%	17/83 = 20%	31/83 = 37%
CWO	RoG-Joint	0/2848 = 0%	3651/6709 = 54%	305/3058 = 10%	1293/3651 = 35%	0/683 = 0%	1835/2729 = 67%	363/2729 = 13%
	RoG-Sep	0/2848 = 0%	3183/6129 = 52%	341/2946 = 12%	1295/3183 = 41%	0/683 = 0%	1840/2620 = 70%	268/2620 = 10%
CQ	SubgraphRAG (Llama3.1-8B)	210/2848 = 7%	2417/5028 = 48%	21/2611 = 1%	98/2417 = 4%	199/683 = 29%	114/1052 = 11%	819/1052 = 78%
	SubgraphRAG (GPT4o-mini)	203/2848 = 7%	2194/5205 = 42%	30/3011 = 1%	159/2194 = 7%	137/683 = 20%	170/953 = 18%	564/953 = 59%

Table 7: Ablation studies with different retrievers, using the same prompt and Llama3.1-8B-Instruct as the reasoner. Rand refers to random triple sampling, RandNoAns removes triples with ground-truth answers after random sampling, and NoRetriever directly asks questions without KG info.

	WebQ	SP	CWO	Q		Web	QSP-sub				CW	/Q-sub		
	Macro-F1	Hit	Macro-F1	Hit	Macro-F1	Micro-F1	Hit	Hit@1	Score _h	Macro-F1	Micro-F1	Hit	Hit@1	Score
SubgraphRAG + Rand	37.69	60.14	27.34	35.85	37.79	17.79	60.74	54.97	65.06	29.97	29.15	39.15	35.11	52.45
SubgraphRAG + RandNoAns	21.18	33.54	16.40	22.71	20.61	8.64	33.03	27.33	47.29	16.47	16.44	23.00	19.42	43.79
SubgraphRAG + NoRetriever	35.86	51.90	25.64	32.34	35.03	17.01	51.38	47.59	55.57	27.42	22.66	33.95	30.65	44.87
SubgraphRAG + cosine similarity	58.41	74.14	34.59	43.61	59.26	37.31	75.43	71.20	73.16	39.05	35.48	49.02	43.68	55.72
SubgraphRAG + Retrieve-Rewrite-Answer	8.96	11.43	-	-	9.11	5.23	11.43	10.84	63.45	-	-	-	-	-
SubgraphRAG + StructGPT	62.14	75.00	-	-	62.55	44.72	75.69	73.57	80.82	-	-	-	-	-
SubgraphRAG + G-Retriever	48.91	64.50	28.47	34.58	49.92	28.08	65.88	62.60	72.54	31.55	32.22	38.17	35.22	58.22
SubgraphRAG + RoG-Sep	57.68	74.39	36.85	45.23	59.36	40.32	76.65	72.61	79.69	44.11	43.69	54.04	47.68	66.18
SubgraphRAG	70.57	86.61	47.16	56.98	72.10	46.56	88.58	84.80	82.42	54.76	51.76	65.80	59.69	62.89

substantially from more powerful LLMs. Our method also excels in truth-grounded QA, with Score_h consistently outperforming baselines by up to 12%, driven by our prompt design that encourages explicit reasoning based on retrieved triples.

Additionally, our framework generalizes well, with retrievers trained on one dataset performing effectively on others. Although moderate performance decay occurs due to domain shift, this can largely be mitigated by including more triples extracted by the retriever. Notably, the decay is generally minor on WebQSP and WebQSP-sub but more pronounced on CWQ and CWQ-sub, potentially due to greater label leakage from the WebQSP test set to the CWQ training set.

In terms of efficiency, studies indicate that overall latency for an LLM to answer a question is primarily dominated by the number of LLM calls, followed by output token count, with input token count having minimal impact ³. SubgraphRAG, despite potentially varying input token counts, utilizes only one LLM call per question, whereas various baselines require multiple calls.

Multi-Hop Performance Breakdown. Table 5 presents the performance breakdown by reasoning hops, with RoG as the primary baseline. Even with an 8B LLM, our method significantly outperforms RoG on multi-hop reasoning questions. This improvement can be attributed to our design approach, which avoids constraints on the types of retrieved subgraphs, allowing the LLMs' reasoning capabilities to be better utilized. For 1-hop questions, SubgraphRAG outperforms the baselines on WebQSP, though its performance is slightly lower on CWQ.

Truth-grounded QA Analysis. Table 6 provides a detailed analysis of truth-grounding in generated answers, showing that previous methods often produce correct answers that are unsupported by the retrieved results, increasing the risk of hallucination. In particular, our method is significantly less likely to generate answers that are not present in the retriever results (the NR answers in the table). For instance, on CWQ-sub, more than 10% of RoG's correct answers are NR, while SubgraphRAG keeps this to just 1%. SubgraphRAG can also decline to answer when there is insufficient evidence, with refusal rates increasing from 2% to 19% on WebQSP and from 7% to 29% on CWQ for questions lacking answers in the KG, compared to questions with KG-supported answers. In contrast, baseline models consistently provide answers even when supporting evidence is absent. Together, these qualities make SubgraphRAG a more trustworthy and truth-grounded KGQA framework.

Explainability. By retrieving flexible, high-quality evidence subgraphs as context and exploiting the strong reasoning capabilities of pre-trained LLMs, SubgraphRAG can natively provide effective explanations along with answer predictions. In contrast, explainable predictions with fine-tuned

 $^{^3}$ For LLMs like GPT-4 and Claude-3.5, an additional input token typically adds about $10^{-2} \times$ the latency of an extra output token and $10^{-4} \times$ that of an additional LLM call (Vivek, 2024).

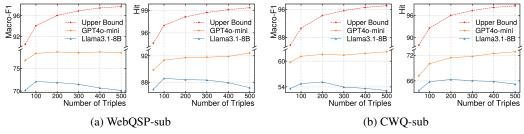


Figure 4: Ablation studies on the number of retrieved triples used in LLM reasoners.

LLMs necessitate extra labeling efforts for preserving explainability (Luo et al., 2024) or post hoc explainability approaches (Krishna et al., 2023). To illustrate the explanability of our framework, we provide multiple examples of our reasoner's responses (see Appendix E).

Performance with Different # of Retrieved Triples. It has been observed that irrelevant context can mislead LLMs (Wu et al., 2024), suggesting that retrieving more triples may not always improve performance. To investigate this, we conducted experiments (Fig. 4) using Llama3.1-8B-Instruct and GPT-40-mini with varying numbers of retrieved triples. The upper bound is defined such that if the answer entity is present in the retrieved results, it is considered a correct answer. The results show that for the less powerful LLM, Llama3.1-8B-Instruct, performance initially improves with more triples but then declines, indicating the model is misled by irrelevant information and struggles to extract key triples from larger contexts. In contrast, GPT-40-mini's performance consistently improves with more triples, demonstrating greater robustness against irrelevant context. These findings suggest that LLMs have varying capacities for handling retrieved information, and the number of retrieved triples should be adjusted to align with the model's capacity and available computational resources.

Performance with Different Retrievers. To assess the impact of our retriever on the overall performance, we conducted extensive studies by replacing the retrieval results with those from baseline retrievers, while keeping all other settings constant (same prompts, same LLM reasoners). Tables 7 and 9 show the results using Llama3.1-8B-Instruct and GPT-4o-mini as the LLM reasoners, respectively. For retrievers from Retrieve-Rewrite-Answer, StructGPT, G-Retriever, and RoG, we used all their retrieved triples, while for other cases, we kept 100 triples. Clearly, pairing with our retriever yields significantly better performance on all metrics and datasets compared to using baseline retrievers, indicating that our superior results are largely due to the effectiveness of our improved retriever.

5 CONCLUSION

This paper presents SubgraphRAG, a novel framework for KG-based RAG. It performs efficient and flexible subgraph retrieval followed by prompting unfine-tuned LLMs for reasoning. SubgraphRAG demonstrates better or comparable accuracy, efficiency, and explainability compared to existing KG-based RAG approaches.

ACKNOWLEDGEMENT

M. Li, S. Miao, and P. Li are partially supported by NSF awards PHY-2117997, IIS-2239565, IIS-2428777, and CCF-2402816; DOE award DE-FOA-0002785; JPMC faculty awards; OpenAI Researcher Access Program Credits; and Microsoft Azure Research Credits for Generative AI.

REFERENCES

Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, pp. 1247–1250, 2008.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego De Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol

- Vinyals, Simon Osindero, Karen Simonyan, Jack Rae, Erich Elsen, and Laurent Sifre. Improving language models by retrieving from trillions of tokens. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 2206–2240, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, pp. 1877–1901, 2020.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4. arXiv preprint arXiv:2303.12712, 2023.
- Michel Chein and Marie-Laure Mugnier. *Graph-based knowledge representation: computational foundations of conceptual graphs.* Springer Science & Business Media, 2008.
- Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count substructures? In *Advances in Neural Information Processing Systems*, pp. 10383–10395, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina 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)*, pp. 4171–4186, 2019.
- Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273, 2022.
- Abhimanyu Dubey et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- Yifu Gao, Linbo Qiao, Zhigang Kan, Zhihua Wen, Yongquan He, and Dongsheng Li. Two-stage generative question answering on temporal knowledge graph using large language models. *arXiv* preprint arXiv:2402.16568, 2024a.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2024b.
- Tiezheng Guo, Qingwen Yang, Chen Wang, Yanyi Liu, Pan Li, Jiawei Tang, Dapeng Li, and Yingyou Wen. Knowledgenavigator: Leveraging large language models for enhanced reasoning over knowledge graph. *Complex & Intelligent Systems*, pp. 1–14, 2024.
- Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using networkx. In *Proceedings of the 7th Python in Science Conference*, pp. 11–15, 2008.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, 2017.
- Taher H. Haveliwala. Topic-sensitive pagerank. In *Proceedings of the 11th International Conference on World Wide Web*, pp. 517–526, 2002.

- Gaole He, Yunshi Lan, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. Improving multi-hop knowledge base question answering by learning intermediate supervision signals. In *Proceedings of the 14th ACM international conference on web search and data mining*, pp. 553–561, 2021.
- Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *arXiv* preprint arXiv:2402.07630, 2024.
- Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. Grag: Graph retrieval-augmented generation. *arXiv preprint arXiv:2405.16506*, 2024.
- Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 1049–1065, 2023.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv* preprint *arXiv*:2311.05232, 2023.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 2023.
- Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong 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, pp. 9237–9251, 2023a.
- Jinhao Jiang, Kun Zhou, Xin Zhao, and Ji-Rong Wen. UniKGQA: Unified retrieval and reasoning for solving multi-hop question answering over knowledge graph. In *International Conference on Learning Representations*, 2023b.
- Jinhao Jiang, Kun Zhou, Wayne Xin Zhao, Yang Song, Chen Zhu, Hengshu Zhu, and Ji-Rong Wen. Kg-agent: An efficient autonomous agent framework for complex reasoning over knowledge graph. *arXiv preprint arXiv:2402.11163*, 2024.
- Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Suhang Wang, Yu Meng, and Jiawei Han. Graph chain-of-thought: Augmenting large language models by reasoning on graphs. *arXiv preprint arXiv:2404.07103*, 2024.
- Minki Kang, Jin Myung Kwak, Jinheon Baek, and Sung Ju Hwang. Knowledge graph-augmented language models for knowledge-grounded dialogue generation. *arXiv preprint arXiv:2305.18846*, 2023.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6769–6781, 2020.
- Jungo Kasai, Keisuke Sakaguchi, yoichi takahashi, Ronan Le Bras, Akari Asai, Xinyan Velocity Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. Realtime QA: What's the answer right now? In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi. Kg-gpt: A general framework for reasoning on knowledge graphs using large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, 2022.

- Satyapriya Krishna, Jiaqi Ma, Dylan Z Slack, Asma Ghandeharioun, Sameer Singh, and Himabindu Lakkaraju. Post hoc explanations of language models can improve language models. In *Advances in Neural Information Processing Systems*, 2023.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano, Sean Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, Daniel Haziza, Luca Wehrstedt, Jeremy Reizenstein, and Grigory Sizov. xformers: A modular and hackable transformer modelling library. https://github.com/facebookresearch/xformers, 2022.
- Quinn Leng, Jacob Portes, Sam Havens, Matei Zaharia, and Michael Carbin. Long context rag performance of large language models. *arXiv preprint arXiv:2411.03538*, 2024.
- Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. Distance encoding: Design provably more powerful neural networks for graph representation learning. *Advances in Neural Information Processing Systems*, 33, 2020.
- Shiyang Li, Yifan Gao, Haoming Jiang, Qingyu Yin, Zheng Li, Xifeng Yan, Chao Zhang, and Bing Yin. Graph reasoning for question answering with triplet retrieval. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 3366–3375, 2023a.
- Xianzhi Li, Samuel Chan, Xiaodan Zhu, Yulong Pei, Zhiqiang Ma, Xiaomo Liu, and Sameena Shah. Are ChatGPT and GPT-4 general-purpose solvers for financial text analytics? a study on several typical tasks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 408–422, 2023b.
- Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*, 2023c.
- Zijian Li, Qingyan Guo, Jiawei Shao, Lei Song, Jiang Bian, Jun Zhang, and Rui Wang. Graph neural network enhanced retrieval for question answering of llms. *arXiv preprint arXiv:2406.06572*, 2024.
- Ke Liang, Lingyuan Meng, Meng Liu, Yue Liu, Wenxuan Tu, Siwei Wang, Sihang Zhou, Xinwang Liu, Fuchun Sun, and Kunlun He. A survey of knowledge graph reasoning on graph types: Static, dynamic, and multi-modal. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. Kagnet: Knowledge-aware graph networks for commonsense reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2829–2839, 2019.
- Guangyi Liu, Yongqi Zhang, Yong Li, and Quanming Yao. Explore then determine: A gnn-llm synergy framework for reasoning over knowledge graph. arXiv preprint arXiv:2406.01145, 2024a.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024b.
- Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *International Conference on Learning Repre*sentations, 2024.
- Shengjie Ma, Chengjin Xu, Xuhui Jiang, Muzhi Li, Huaren Qu, and Jian Guo. Think-on-graph 2.0: Deep and interpretable large language model reasoning with knowledge graph-guided retrieval. *arXiv preprint arXiv:2407.10805*, 2024.

- Costas Mavromatis and George Karypis. Gnn-rag: Graph neural retrieval for large language model reasoning. *arXiv preprint arXiv:2405.20139*, 2024.
- Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: higher-order graph neural networks. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, 2019.
- Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. MTEB: Massive text embedding benchmark. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 2014–2037, 2023.
- Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*, 36(7):3580–3599, 2024.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pp. 8024–8035, 2019.
- Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. Graph retrieval-augmented generation: A survey. *arXiv preprint arXiv:2408.08921*, 2024.
- Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, 2009.
- Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. Okapi at trec-3. In *TREC*, volume 500-225 of *NIST Special Publication*, pp. 109–126, 1994.
- Ian Robinson, Jim Webber, and Emil Eifrem. *Graph databases: new opportunities for connected data.* "O'Reilly Media, Inc.", 2015.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmentation reduces hallucination in conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 3784–3803, 2021.
- Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William Cohen. Open domain question answering using early fusion of knowledge bases and text. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4231–4242, 2018.
- Haitian Sun, Tania Bedrax-Weiss, and William Cohen. Pullnet: Open domain question answering with iterative retrieval on knowledge bases and text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2380–2390, 2019.
- Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel Ni, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph. In *International Conference on Learning Representations*, 2024a.
- Lei Sun, Zhengwei Tao, Youdi Li, and Hiroshi Arakawa. Oda: Observation-driven agent for integrating llms and knowledge graphs. *arXiv preprint arXiv:2404.07677*, 2024b.
- Alon Talmor and Jonathan Berant. The web as a knowledge-base for answering complex questions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 641–651, 2018.

- Dhaval Taunk, Lakshya Khanna, Siri Venkata Pavan Kumar Kandru, Vasudeva Varma, Charu Sharma, and Makarand Tapaswi. Grapeqa: Graph augmentation and pruning to enhance question-answering. In Companion Proceedings of the ACM Web Conference 2023, pp. 1138–1144, 2023.
- Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-evolve: Deep temporal reasoning for dynamic knowledge graphs. In *international conference on machine learning*, pp. 3462–3471. PMLR, 2017.
- Vivek. Understanding input/output tokens vs latency tradeoff. https://minusx.ai/blog/input-vs-output-tokens/, 2024. URL https://minusx.ai/blog/input-vs-output-tokens/. Accessed: 2024-10-27.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. FreshLLMs: Refreshing large language models with search engine augmentation. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 13697–13720, 2024.
- Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering. *arXiv preprint arXiv:2308.13259*, 2023a.
- Xintao Wang, Qianwen Yang, Yongting Qiu, Jiaqing Liang, Qianyu He, Zhouhong Gu, Yanghua Xiao, and Wei Wang. Knowledgpt: Enhancing large language models with retrieval and storage access on knowledge bases. *arXiv preprint arXiv:2308.11761*, 2023b.
- Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. Knowledge graph prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19206–19214, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems, 2022.
- Yilin Wen, Zifeng Wang, and Jimeng Sun. Mindmap: Knowledge graph prompting sparks graph of thoughts in large language models. *arXiv preprint arXiv:2308.09729*, 2023.
- Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. Transfertransfo: A transfer learning approach for neural network based conversational agents. *arXiv* preprint arXiv:1901.08149, 2019.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, 2020.
- Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai Zhang, and Yanghua Xiao. How easily do irrelevant inputs skew the responses of large language models? *arXiv preprint arXiv:2404.03302*, 2024.
- Yike Wu, Nan Hu, Guilin Qi, Sheng Bi, Jie Ren, Anhuan Xie, and Wei Song. Retrieve-rewrite-answer: A kg-to-text enhanced llms framework for knowledge graph question answering. *arXiv* preprint arXiv:2309.11206, 2023.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee, Chen Zhu, Zihan Liu, Sandeep Subramanian, Evelina Bakhturina, Mohammad Shoeybi, and Bryan Catanzaro. Retrieval meets long context large language models. In *International Conference on Learning Representations*, 2024.

- Xiao Yang, Kai Sun, Hao Xin, Yushi Sun, Nikita Bhalla, Xiangsen Chen, Sajal Choudhary, Rongze Daniel Gui, Ziran Will Jiang, Ziyu Jiang, et al. Crag-comprehensive rag benchmark. arXiv preprint arXiv:2406.04744, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems*, 2023.
- Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. Qa-gnn: Reasoning with language models and knowledge graphs for question answering. In *North American Chapter of the Association for Computational Linguistics (NAACL)*, 2021.
- Wen-tau Yih, Matthew Richardson, Chris Meek, Ming-Wei Chang, and Jina Suh. The value of semantic parse labeling for knowledge base question answering. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 201–206, 2016.
- Jing Zhang, Xiaokang Zhang, Jifan Yu, Jian Tang, Jie Tang, Cuiping Li, and Hong Chen. Subgraph retrieval enhanced model for multi-hop knowledge base question answering. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5773–5784, 2022.
- Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. Labeling trick: A theory of using graph neural networks for multi-node representation learning. *Advances in Neural Information Processing Systems*, 34:9061–9073, 2021.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren's song in the ai ocean: A survey on hallucination in large language models. *arXiv* preprint arXiv:2309.01219, 2023.

A ADDITIONAL RELATED WORK

The field of KGQA has evolved significantly over time. Early approaches to KGQA do not rely on LLMs for answer generation (Yasunaga et al., 2021; Taunk et al., 2023; Zhang et al., 2022; Lin et al., 2019; Sun et al., 2019), though they often employ pre-trained language models (PLMs) like BERT as text encoders (Devlin et al., 2019). These methods typically search for answer entities with GNNs Yasunaga et al. (2021); Taunk et al. (2023); Lin et al. (2019) or LSTM-based models Sun et al. (2019); Zhang et al. (2022).

With the rapid advancement of LLMs in recent years, researchers began to leverage them in KGQA. For instance, recent work has explored using LLMs as translators, converting natural language questions into executable SQL queries for KG databases to retrieve the answers (Jiang et al., 2023a; Wang et al., 2023b).

Contemporary approaches have further expanded the role of LLMs, utilizing them for both knowledge retrieval from KGs and reasoning (Kim et al.; Gao et al., 2024a; Wang et al., 2024; Guo et al., 2024; Ma et al., 2024; Sun et al., 2024a; Jiang et al., 2024; Jin et al., 2024). The strength of this strategy lies in LLMs' ability to handle multi-hop tasks by breaking them down into manageable steps. However, as discussed in Section 1, this often necessitates multiple LLM calls, resulting in high latency. To mitigate this issue, some frameworks have attempted to fine-tune LLMs to memorize knowledge, but this reduces their ability to generalize to dynamically updated or novel KGs (Luo et al., 2024; Mavromatis & Karypis, 2024). Other models have explored fine-tuning adapters embedded in fixed LLMs to better preserve their general reasoning capabilities while adapting to specific KGs (He et al., 2024; Gao et al., 2024a; Hu et al., 2024).

In parallel with these developments, several approaches have emerged that, like our approach, allow LLMs to reason over subgraphs (Kim et al.; Liu et al., 2024a; Li et al., 2023a; Guo et al., 2024; Wu et al., 2023; Li et al., 2024; Wen et al., 2023), though they employ different retrieval strategies. Kim et al. breaks queries into sub-queries, retrieving evidence for each sub-query before reasoning over the collected evidence. Guo et al. (2024) uses PLM-based entity extraction followed by multi-hop expansion for retrieval. Li et al. (2023a) linearizes KG triples into natural language for global triple retrieval using BM25 or dense passage retrievers. Wen et al. (2023) extracts topic entities and merges triples into the retrieved subgraph using two heuristic methods: connecting topic entities via paths and retrieving their 1-hop neighbors.

While these subgraph retrieval strategies share similarities with SubgraphRAG, they lack several key advantages of it, such as retrieval efficiency, adjustable subgraph sizes, and flexible subgraph types. Consequently, they frequently result in suboptimal coverage of relevant information in the retrieved subgraphs. SubgraphRAG addresses these limitations, offering a more comprehensive and adaptable approach to KGQA that builds upon and extends the capabilities of existing methods, while fully leveraging the power of advanced LLMs.

B ADDITIONAL EXPERIMENT DETAILS FOR SUBGRAPH RETRIEVAL

B.1 ADDITIONAL IMPLEMENTATION AND EVALUATION DETAILS

For Retrieve-Rewrite-Answer, RoG, and G-Retriever, we utilize their official open-source implementations for training and evaluation. While a pre-trained RoG model is publicly available, it was jointly trained on the training subset of both WebQSP and CWQ, causing a label leakage issue due to sample duplication across the two datasets. To address this, we retrain the RoG model separately on the training subset of WebQSP and CWQ. Both Retrieve-Rewrite-Answer and G-Retriever were originally only evaluated on the WebQSP dataset. We managed to adapt the G-Retriever codebase for an evaluation on the CWQ dataset.

RoG, G-Retriever, the cosine similarity baseline, and all SubgraphRAG variants perform relevant subgraph retrieval from the identical rough subgraphs. These subgraphs are centered around the topic entities within a maximum number of hops, and they are included in the released RoG implementation. In contrast, Retrieve-Rewrite-Answer directly loads and queries the raw KG with a database server.

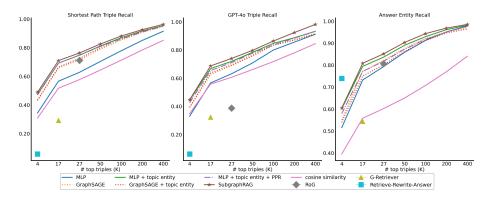


Figure 5: Retrieval effectiveness on WebQSP across a spectrum of K values for top-K triple retrieval.

Table 8: Breakdown of recall evaluation for CWQ over topic entity count. Best results are in **bold**.

Model	Shortest F	ath Triples	GPT-40-la	beled Triples	Answer	Entities
	Single (46.3%)	Multiple (53.7%)	Single (46.3%)	Multiple (53.7%)	Single (46.3%)	Multiple (53.7%)
cosine similarity	0.586	0.403	0.587	0.499	0.517	0.638
RoG	0.711	0.547	0.324	0.275	0.769	0.903
G-Retriever	0.245	0.128	0.289	0.156	0.393	0.360
MLP	0.788	0.568	0.777	0.600	0.894	0.873
MLP + topic entity	0.858	0.690	0.819	0.718	0.894	0.889
SubgraphRAG	0.877	0.754	0.871	0.820	0.911	0.916

GraphSAGE originally only deals with node attributes and we propose a straightforward extension of it for handling both the node attributes and edge attributes. Let \mathbf{z}_e be the text embedding of an entity $e \in \mathcal{E}$ and \mathbf{z}_r be the text embedding of a relation $r \in \mathcal{R}$. A GraphSAGE layer updates entity representations with

$$\mathbf{z}_{\mathcal{N}(e)}^{l+1} = \mathbf{MEAN}\left(\left\{ \left[\mathbf{z}_{e'}^{l} \middle| \mathbf{z}_{r}\right] \middle| \left(e', r, e\right) \in \mathcal{G}\right\}\right)$$
(3)

$$\mathbf{z}_{e}^{l+1} = \sigma\left(\left[\mathbf{z}_{e}^{l} \| \mathbf{z}_{\mathcal{N}(e)}^{l+1}\right)$$

$$\mathbf{z}_{e}^{0} = \mathbf{z}_{e}$$

$$(4)$$

$$\mathbf{z}_e^0 = \mathbf{z}_e \tag{5}$$

where $\sigma(\cdot)$ is an MLP. Empirically we find that 1-layer GraphSAGE yields the best performance.

B.2 IMPLEMENTATION DEPENDENCIES

Our implementation is based on the following packages: PyTorch (Paszke et al., 2019), Transformers (Wolf et al., 2020), xFormers (Lefaudeux et al., 2022), NetworkX (Hagberg et al., 2008), and PyTorch Geometric (Fey & Lenssen, 2019). We employ the built-in implementation of PPR from NetworkX.

B.3 ADDITIONAL EXPERIMENT RESULTS FOR VARYING K IN TOP-K TRIPLE RETRIEVAL

Fig. 5 presents the performance of SubgraphRAG variants on WebQSP across a spectrum of K values for top-K triple retrieval, which is consistent with the observations made for CWQ.

EVALUATION BREAKDOWN OVER TOPIC ENTITY COUNT

See table 8.

Table 9: Ablation studies with different retrievers, using the same prompt and GPT-4o-mini as the reasoner. Rand refers to random triple sampling, RandNoAns removes triples with ground-truth answers after random sampling, and NoRetriever directly asks questions without KG info.

	WebQ	SP	CWO	2		Web	QSP-sub				CW	/Q-sub		
	Macro-F1	Hit	Macro-F1	Hit	Macro-F1	Micro-F1	Hit	Hit@1	Score _h	Macro-F1	Micro-F1	Hit	Hit@1	Score _h
SubgraphRAG + Rand	47.70	68.37	33.13	39.82	47.15	23.74	68.57	64.59	69.67	35.55	34.26	42.94	40.48	54.81
SubgraphRAG + RandNoAns	36.83	49.63	25.69	30.70	35.77	14.62	48.94	44.96	57.12	26.25	26.11	31.60	29.21	48.80
SubgraphRAG + NoRetriever	47.49	71.01	33.43	42.25	46.68	25.53	70.94	62.67	60.22	35.66	31.11	44.66	39.99	44.24
SubgraphRAG + cosine similarity	64.94	78.19	41.23	49.22	65.26	43.84	78.90	74.15	73.03	45.07	42.27	53.83	49.37	57.06
SubgraphRAG + Retrieve-Rewrite-Answer	38.26	53.62	-	-	37.54	19.62	53.37	50.10	67.68	-	-	-	-	-
SubgraphRAG + StructGPT	71.62	82.68	-	-	71.66	59.05	82.87	80.44	81.48	-	-	-	-	-
SubgraphRAG + G-Retriever	59.28	75.25	36.79	42.23	59.55	35.19	76.01	72.87	76.31	39.55	39.42	45.47	43.29	57.67
SubgraphRAG + RoG-Sep	70.08	82.25	44.30	51.15	70.91	54.69	83.45	78.51	81.00	49.85	50.40	57.51	52.35	64.27
SubgraphRAG	77.45	90.11	54.13	62.02	78.34	58.44	91.34	87.36	82.21	61.13	58.86	70.01	65.48	64.20

C SUPPLEMENTARY KGQA RESULTS

To obtain a single metric to assess the degree of truth-grounded QA performance, we follow (Yang et al., 2024) to develop a scoring strategy to aggregate the results of these tables. Specifically, for samples with answer entities present in the KG, correct answers receive a score of +1, incorrect answers -1, and missing answers 0. For samples without answer entities in the KG, missing answers, answers with entities found in the retrieved results, and answers with entities not found in the retrieved results are scored +1, -1, and -1.5, respectively. Denote s_i as the score yielded for sample i and a_i the number of answers outputted by the LLM reasoner for this sample. For a dataset with n samples: $\operatorname{Score}_h = \operatorname{Normalize}\left(\frac{1}{n}\sum_{i=0}^n \frac{s_i}{a_i}\right)$, where we use min-max normalization with min = 0 and max = 100. This metric penalizes hallucinated answers while favoring missing answers over incorrect ones.

To assess the impact of our retriever on overall performance as analyzed in Sec. 4.1, Table 9 provides additional ablation results using GPT-4o-mini as the reasoner.

D PROMPT TEMPLATES

Fig. 6 is the detailed prompt template used in our experiments, and Fig. 7 provides the prompt employed to label relevant triples via GPT-4o.

```
Input Prompts for KGQA
Based on the triples retrieved from a knowledge graph, please answer the question. Please return formatted answers as a list, each
prefixed with "ans:"
User: // ICL example
Triplets:
(Lou Seal.sports.mascot.team.San Francisco Giants)
(San Francisco Giants, sports_sports_team.championships, 2012 World Series)
(San Francisco Giants, sports_championship_event.champion, 2014 World Series)
(San Francisco Giants,time.participant.event,2014 Major League Baseball season)
(San Francisco Giants,time.participant.event,2010 World Series)
(San Francisco Giants, time.participant.event, 2010 Major League Baseball season)
(San Francisco Giants, sports_sports_team.championships, 2014 World Series)
(San Francisco Giants, sports_team.team_mascot, Crazy Crab)
(San\ Francisco\ Giants, sports\_team.championships, 2010\ World\ Series)
(San\ Francisco\ Giants, sports.professional\_sports\_team.owner\_s, Bill\ Neukom)
(San Francisco Giants, time.participant.event, 2012 World Series)
(San Francisco, sports_sports_team_location.teams, San Francisco Giants)
(San Francisco Giants, sports_team.arena_stadium, AT&T Park)
(AT&T Park,location.location.events,2012 World Series)
(m.011zsc4_,organization.leadership.organization,San Francisco Giants)
(San Francisco Giants, sports_team.previously_known_as, New York Giants)
(AT&T Park,location.location.events,2010 World Series)
(Crazy Crab, sports.mascot.team, San Francisco Giants)
(New York Giants,baseball.baseball_team.league,National League)
(San Francisco Giants, sports_team.colors, Black)
(San\ Francisco\ Giants, sports\_team.previously\_known\_as, New\ York\ Gothams)
(m.0k079qm,base.schemastaging.team_training_ground_relationship.team,San Francisco Giants)
(m.0k079ry, base.schemastaging.team\_training\_ground\_relationship.team, San\ Francisco\ Giants)
(2010 World Series, time.event.locations, AT&T Park)
(San Francisco Giants,time.participant.event,2012 Major League Baseball season)
(San Francisco Giants,baseball.baseball_team.league,National League)
(m.0crtd80,sports.sports_league_participation.league,National League West)
(San Francisco Giants, sports_team.location, San Francisco)
(San Francisco Giants, sports_team.sport, Baseball)
(m.05n6dtn,baseball.baseball_team_stats.team,San Francisco Giants)
What year did the team with mascot named Lou Seal win the World Series?
Assistant: // ICL example
To find the year the team with mascot named Lou Seal won the World Series, we need to find the team with mascot named Lou
Seal and then find the year they won the World Series.
From the triplets, we can see that Lou Seal is the mascot of the San Francisco Giants.
Now, we need to find the year the San Francisco Giants won the World Series.
From the triplets, we can see that San Francisco Giants won the 2010 World Series and 2012 World Series and 2014 World Series.
So, the team with mascot named Lou Seal (San Francisco Giants) won the World Series in 2010, 2012, and 2014.
Therefore, the formatted answers are:
ans: 2014 (2014 World Series)
ans: 2012 (2012 World Series)
ans: 2010 (2010 World Series)
User: // Sample to infer
Triplets:
(e_a, r_{ab}, e_b) // Retrieved triples
(\mathbf{e}_c, \mathbf{r}_{cd}, \mathbf{e}_d)
Ouestion:
what's ...? // Q in datasets
```

Figure 6: Detailed prompt for LLM question-answering

Input Prompts for Labeling Relevant Triples via GPT-40 System: Based on the triplets retrieved from a knowledge graph, please select relevant triplets for answering the question. Please return formatted triplets as a list, each prefixed with "evidence:". User: // Sample to infer Triplets: (e_a, r_{ab}, e_b) // Retrieved triples (e_c, r_{cd}, e_d) ... Question: what's ...? // Q in datasets

Figure 7: Detailed prompt for labeling triples via GPT-4o.

E EXPLAINABLITY EXAMPLES

To demonstrate the superior explainability of our approach, we provide multiple example responses from our LLM reasoners below, which cover questions requiring different reasoning hops and logic chains. We also have included examples where our LLM reasoners refused to answer due to insufficient information, which helps to provide more truth-grounded answers.

CWQ-WebQTest-100_bf66fd89b6e9fc5fcb96c8b3f7a0e616

Ouestion:

The people from the country that contains Nord-Ouest Department speak what languages today?

LLM Response:

To find the languages spoken in the country that contains Nord-Ouest Department, we need to find the country that contains Nord-Ouest Department.

From the triplets, we can see that Nord-Ouest Department is contained by Haiti.

So, the country that contains Nord-Ouest Department is Haiti.

From the triplets, we can see that Haiti speaks Haitian Creole and French.

Therefore, the people from the country that contains Nord-Ouest Department (Haiti) speak Haitian Creole and French today. The formatted answer is:

ans: Haitian Creole

ans: French

Ground-truth:

Haitian Creole

French

CWQ-WebQTrn-2152_92fba37c9723caee68665ad9a5e4a468

Question:

Which of team owners Tom Hicks sports team are in the American League West?

LLM Response:

To find the teams owned by Tom Hicks that are in the American League West, we need to find the teams owned by Tom Hicks and then find the teams in the American League West.

From the triplets, we can see that Tom Hicks owns the Texas Rangers and Dallas Stars.

Now, we need to find the teams in the American League West

From the triplets, we can see that the teams in the American League West are Texas Rangers, Los Angeles Angels of Anaheim, Seattle Mariners, Oakland Athletics, and Seattle Pilots.

So, among the teams owned by Tom Hicks (Texas Rangers and Dallas Stars), only Texas Rangers is in the American League West. Therefore, the formatted answer is:

ans: Texas Rangers

Ground-truth:

Texas Rangers

CWQ-WebQTrn-557_960c16ffdb29e173df0577fc76c7455d

Question

In the film with the character named Winged Monkey #7, who played Dorothy?

LLM Response:

To find the actor who played Dorothy in the film with the character named Winged Monkey #7, we need to find the film that Winged Monkey #7 is in and then find the actor who played Dorothy in that film.

From the triplets, we can see that Winged Monkey #7 is in the film "The Wizard of Oz" (m.0yj6tr7).

Now, we need to find the actor who played Dorothy in "The Wizard of Oz".

From the triplets, we can see that Judy Garland (m.0k6y72) played Dorothy in "The Wizard of Oz".

Therefore, the formatted answer is:

ans: Judy Garland

Ground-truth:

Judy Garland

CWQ-WebQTest-1320_c5498ca807d2e1ec30d4c8fdd41f0bf7

Ouestion:

What stadium does the team with mascot named Hank play at?

LLM Response:

To find the stadium where the team with mascot named Hank plays, we need to find the team with mascot named Hank and then find the stadium they play at.

From the triplets, we can see that Hank is the mascot of the Milwaukee Brewers.

Now, we need to find the stadium where the Milwaukee Brewers play.

From the triplets, we can see that the Milwaukee Brewers play at Miller Park.

Therefore, the formatted answer is:

ans: Miller Park

Ground-truth:

Miller Park

CWQ-WebQTrn-2904_cb8d941f825c5dd7b2d6f3faef7f9229

Ouestion

What city is home for Lauren Conrad and had a mayor named Willie Brown?

LLM Response:

To find the city that is home for Lauren Conrad and had a mayor named Willie Brown, we need to find the city where Lauren Conrad lives and the city where Willie Brown was the mayor.

From the triplets, we can see that Lauren Conrad lived in Los Angeles (m.043jjv4) and San Francisco (m.04hy_14).

From the triplets, we can see that Willie Brown was the mayor of San Francisco.

Therefore, the city that is home for Lauren Conrad and had a mayor named Willie Brown is San Francisco.

So, the formatted answer is:

ans: San Francisco

Ground-truth:

San Francisco

CWQ-WebQTrn-124_405a78f132e69f4771aff40d07cd0294

Question:

Which movie did Angelina Jolie direct, that was produced by Erwin Stoff?

LLM Response:

To find the movie that Angelina Jolie directed, which was produced by Erwin Stoff, we need to find the movies that Angelina Jolie directed and the movies that Erwin Stoff produced.

From the triplets, we can see that Angelina Jolie directed the following movies:

- In the Land of Blood and Honey
- Unbroken
- By the Sea

From the triplets, we can see that Erwin Stoff produced the following movies:

- Unbroken
- The Day the Earth Stood Still
- Edge of Tomorrow
- Water for Elephants
- Austin Powers: The Spy Who Shagged Me
- The Blind Side
- The Replacements
- I Am Legend
- Hardball
- The Lake House
- Street Kings
- 47 Ronin

Now, we need to find the common movie between Angelina Jolie's directed movies and Erwin Stoff's produced movies.

The common movie is Unbroken.

Therefore, the formatted answer is:

ans: Unbroken

Ground-truth:

Unbroken

WebQSP-WebQTest-33

Ouestion:

what timezone is sweden

LLM Response:

To find the timezone of Sweden, we can look for the triplets that contain Sweden and have a property related to timezones. From the triplets, we can see that Sweden is located in the Central European Time Zone.

Therefore, the answer is:

ans: Central European Time Zone

Ground-truth:

Central European Time Zone

WebQSP-WebQTest-62

Question:

what was jesse james killed with

LLM Response:

To find what Jesse James was killed with, we need to find the cause of his death.

From the triplets, we can see that Jesse James was killed by Assassination.

Now, we need to find the cause of Assassination.

From the triplets, we can see that Assassination includes causes of death such as Assassination in ways which appear natural,

Assassination by 'Accident', Assassination by firearm, and Assassination by 'Suicide'. However, we are looking for what Jesse James was killed with, not how he was killed.

From the triplets, we can see that Firearm is a cause of death that includes causes of death such as Assassination by firearm.

Therefore, Jesse James was killed with a Firearm.

So, the answer is: ans: Firearm

Ground-truth:

Firearm

WebQSP-WebQTest-188

Question:

where did bristol palin go to school

To find where Bristol Palin went to school, we need to find the educational institutions that Bristol Palin attended.

From the triplets, we can see that Bristol Palin attended the following educational institutions:

- Wasilla High School (m.0byl4zl)
- Juneau-Douglas High School (m.0cpkbwp)
- West Anchorage High School (m.0cpk4bl)

Therefore, the formatted answers are:

ans: Wasilla High School

ans: Juneau-Douglas High School ans: West Anchorage High School

Ground-truth:

West Anchorage High School

Wasilla High School

Juneau-Douglas High School

WebQSP-WebQTest-760

Ouestion:

what party does barack obama belong to

LLM Response:

From the triplets, we can see that Barack Obama is associated with the Democratic Party in several instances:

- 1. (Barack Obama,government.politician.party,m.03gjhww)
- 2. (Barack Obama 2004 US Senate Campaign, government.election_campaign.party, Democratic Party)
- 3. (Barack Obama 1996 Illinois Senate Campaign,government.election_campaign.party,Democratic Party)
- 4. (Barack Obama presidential primary campaign, 2008, government.election_campaign.party, Democratic Party)
- 5. (Democratic Party presidential primaries, 2008,government.primary_election.party,Democratic Party)
- $6. \ (m.03gjhww,government.political_party_tenure.party,Democratic\ Party)$

Therefore, the formatted answer is:

ans: Democratic Party

Ground-truth:

Democratic Party

WebQSP-WebQTest-386

Onestion:

how deep is lake merritt oakland

LLM Response:

To find the depth of Lake Merritt in Oakland, we need to look for any information related to the depth of Lake Merritt.

From the triplets, we can see that Lake Merritt is a lake, but there is no direct information about its depth.

However, we can see that Lake Merritt is a National Historic Landmark (m.06_44f4,

 $protected_sites.natural_or_cultural_site_listing_designation_as_national_or_cultural_site, m.06_44f4) \ and \ a \ National \ Historical \ Additional \ Addition$

Landmark (m.06_44f4, protected_sites.natural_or_cultural_site_listing.category_or_criteria, National Historic Landmark).

We can also see that Lake Merritt is a lake (Lake Merritt, common.topic.notable_types, Lake) and it is located in Oakland (Lake Merritt, location.location.containedby, Oakland).

But there is no information about the depth of Lake Merritt.

However, we can see that there is a mention of a location with geolocation m.0cs7t50 (Lake Merritt, location.location.geolocation, m.0cs7t50), but this does not provide any information about the depth of the lake.

Therefore, based on the given triplets, we cannot find any information about the depth of Lake Merritt.

ans: not available

Ground-truth:

3.048

WebQSP-WebQTest-612

Ouestion:

who does jeremy shockey play for in 2012

LLM Response:

To find out who Jeremy Shockey played for in 2012, we need to find the team he played for and the year he played for that team. From the triplets, we can see that Jeremy Shockey played for the New Orleans Saints, Carolina Panthers, Miami Hurricanes football, and New York Giants.

However, we need to find the year he played for each of these teams.

From the triplets, we can see that Jeremy Shockey played for the New Orleans Saints in 2007, for the Carolina Panthers in 2007,

for the Miami Hurricanes football in 2001,

and for the New York Giants in 2002.

However, there is no information about Jeremy Shockey playing for any of these teams in 2012.

Therefore, we cannot determine which team Jeremy Shockey played for in 2012 based on the given triplets. ans: not available

Ground-truth:

Carolina Panthers

CWQ-WebQTest-397_0b3a8b0643fac530a078012a025201fa

Ouestion:

Harvard Board of Overseers is part of the school with what colors?

LLM Response:

To find the school with the colors that Harvard Board of Overseers is part of, we need to find the parent organization of Harvard Board of Overseers.

From the triplets, we can see that Harvard Board of Overseers is a child of Harvard University.

Now, we need to find the colors of Harvard University.

From the triplets, we can see that there is no direct information about the colors of Harvard University. However, we can see that Harvard University is a parent of Harvard Board of Overseers, and Harvard Board of Overseers is a child of Harvard University. But we can also see that Harvard University is a parent of itself (Harvard

University, education.educational_institution.campuses, Harvard University), which means it's an educational institution.

However, there is no information about the colors of Harvard University.

Ground-truth:

Crimson

CWQ-WebQTest-1528_2f6fb6d585b98261a3ae6d2a112a4c91

Ouestion:

What movie starring Logan Lerman that was edited by Thomas J Nordberg?

LLM Response:

To find the movie starring Logan Lerman that was edited by Thomas J Nordberg, we need to find the movies edited by Thomas J Nordberg and then find the ones starring Logan Lerman.

From the triplets, we can see that the following movies were edited by Thomas J Nordberg:

- Scary Movie 2
- The Host
- I Love You Phillip Morris
- U Turn
- Any Given Sunday
- The Guardian
- Alexander

Now, we need to find the movies among these that star Logan Lerman.

From the triplets, we can see that Logan Lerman starred in the following movies:

- What Women Want
- The Perks of Being a Wallflower
- Fury
- Percy Jackson: Sea of Monsters
- Noah
- Riding in Cars with Boys
- 3:10 to Yuma

However, none of these movies are among the ones edited by Thomas J Nordberg.

But, we can see that Fury is not edited by Thomas J Nordberg but Fury is related to another film edited by Thomas J Nordberg. The film Fury is related to Fury (2014) which is a film edited by Thomas J Nordberg.

Ground-truth:

What Women Want

CWQ-WebQTest-12_c701ad2b5b8ef3f3ed26dd2ed8703d05

Question

Who held his governmental position from December 12, 1808 and was Governor of Ohio in 2011?

LLM Response:

To find the answer, we need to find a person who held a governmental position from December 12, 1808 and was Governor of Ohio in 2011.

From the triplets, we can see that there is no direct information about a person holding a governmental position from December 12, 1808. However, we can find a person who was Governor of Ohio in 2011.

From the triplets, we can see that Mike DeWine was the Governor of Ohio.

However, there is no information about Mike DeWine holding a governmental position from December 12, 1808.

We can also see that there is no information about any other person holding a governmental position from December 12, 1808. Therefore, we cannot find a person who held a governmental position from December 12, 1808 and was Governor of Ohio in 2011. ans: not available

Ground-truth:

Return J. Meigs, Jr.