

A Query-Aware Multi-Path Knowledge Graph Fusion Approach for Enhancing Retrieval-Augmented Generation in Large Language Models

Qikai Wei^a, Huansheng Ning^{a,*}, Chunlong Han^a and Jianguo Ding^b

^aSchool of Computer and Communication Engineering, University of Science and Technology Beijing, Beijing, China

^bBlekinge Institute of Technology, Karlskrona, Sweden

ARTICLE INFO

Keywords:

Retrieval Augmented Generation
Query-Aware Attention
Reward Model
Knowledge Graph Subgraph Fusion
Large Language Models
Retrieval

ABSTRACT

Retrieval Augmented Generation (RAG) has gradually emerged as a promising paradigm for enhancing the accuracy and factual consistency of content generated by large language models (LLMs). However, existing RAG studies primarily focus on retrieving isolated segments using similarity-based matching methods, while overlooking the intrinsic connections between them. This limitation hampers performance in RAG tasks. To address this, we propose QMKGF, a Query-Aware Multi-Path Knowledge Graph Fusion Approach for Enhancing Retrieval Augmented Generation. First, we design prompt templates and employ general-purpose LLMs to extract entities and relations, thereby generating a knowledge graph (KG) efficiently. Based on the constructed KG, we introduce a multi-path subgraph construction strategy that incorporates one-hop relations, multi-hop relations, and importance-based relations, aiming to improve the semantic relevance between the retrieved documents and the user query. Subsequently, we designed a query-aware attention reward model that scores subgraph triples based on their semantic relevance to the query. Then, we select the highest score subgraph and enrich subgraph with additional triples from other subgraphs that are highly semantically relevant to the query. Finally, the entities, relations, and triples within the updated subgraph are utilised to expand the original query, thereby enhancing its semantic representation and improving the quality of LLMs' generation. We evaluate QMKGF on the SQuAD, IIRC, Culture, HotpotQA, and MuSiQue datasets. On the HotpotQA dataset, our method achieves a ROUGE-1 score of 64.98%, surpassing the BGE-Rerank approach by 9.72 percentage points (from 55.26% to 64.98%). Experimental results demonstrate the effectiveness and superiority of the QMKGF approach.

1. Introduction

In recent years, Large Language Models (LLMs) have demonstrated remarkable performance in the field of Natural Language Processing (NLP), finding widespread application across various artificial intelligence tasks [16, 30, 42, 35]. Despite their increasingly powerful generative capabilities, LLMs often produce outputs that appear plausible but are factually incorrect—a phenomenon commonly referred to as hallucination [7, 39].

To address this issue, Retrieval-Augmented Generation (RAG) has emerged as an effective solution [8, 40]. By incorporating external knowledge sources to provide contextual support for generation, RAG significantly mitigates the hallucination phenomenon commonly observed in LLMs and substantially enhances the factual accuracy and reliability of the generated content [37]. Owing to its ability to improve factual consistency, RAG has been widely applied in various domains such as law [3], medical [31], and tourism [29]. Fig. 1 illustrates a comparison between RAG-based approaches and the standard LLM-only paradigm. The green line indicates the LLM-only setting, where the model generates answers directly from the query without

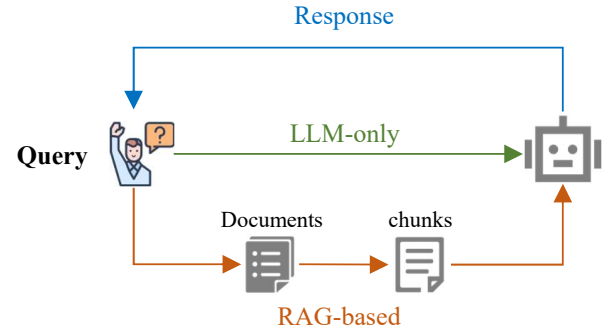


Figure 1: Comparison between LLM-only and RAG-based paradigm.

any external retrieval. The red line represents the RAG-based approach, where a user query is used to retrieve relevant chunks from a document corpus, guiding the LLM to generate more informed responses. The blue line denotes the final response delivered to the user by the LLM.

Existing RAG methods can be divided into document retrieval-based and graph structure-based methods according to the organization of knowledge structure. Document retrieval-based RAG commonly adopt query expansion [27] or re-ranking mechanisms to improve the precision and recall of relevant content. For example, NetEase Youdao [14] proposed the BCEembedding method, which employs a

*Corresponding author

✉ weiqikai@xs.ustb.edu.cn (Q. Wei); ninghuansheng@ustb.edu.cn (H. Ning); hanchunlong@xs.ustb.edu.cn (C. Han); jianguo.ding@bth.se (J. Ding)

ORCID(s): 0000-0002-1048-3814 (Q. Wei); 0000-0001-6413-193X (H. Ning); 0000-0003-1602-1597 (C. Han); 0000-0002-8927-0968 (J. Ding)

two-stage approach: the first stage performs efficient initial retrieval using embeddings, while the second stage applies a RM to conduct fine-grained semantic re-ranking of the retrieved results, thereby enhancing the overall retrieval quality. Xiao et al. [33] introduced the BGE-Rerank model to re-rank the top-k documents retrieved by the embedding model, aiming to enhance the relevance between the retrieved documents and the query. This refinement ultimately boosts the performance of RAG tasks. However, these methods heavily rely on similarity-based query-to-chunk matching, which may introduce irrelevant content and compromise the accuracy of the response. While such approaches are effective for simple tasks, they often fall short when dealing with complex queries [10].

In contrast, graph structure-based retrieval methods for RAG leverage graph augmentation or path-guided mechanisms to enhance the semantic relevance of retrieval results. Representative works include GraphRAG [5], QCG-rerank [28], and KG2RAG [41], which improve content retrieval quality by integrating the semantic relationships between entities. However, existing graph-based RAG approaches typically rely on a single path, which limits the scope of retrieved content and affects the diversity of the final generated output. On the one hand, current methods often construct subgraphs based on fixed rules or structural distances, neglecting the variation in information-bearing capacity across different semantic paths. This makes it difficult to balance semantic diversity and path significance, resulting in the introduction of fragmented and structurally homogeneous information during retrieval. On the other hand, incorporating multiple paths directly may lead to semantic redundancy and the accumulation of noise, which weakens the model's ability to identify critical information. Therefore, if a mechanism can be developed that not only extracts high-quality information from multiple paths but also enables semantic aggregation and filtering, it holds the potential to improve generation quality while mitigating redundancy and noise in the process of graph information incorporation.

To address the above challenges, we propose QMKGF, a Query-aware Multi-path Knowledge Graph Fusion approach. This method enhances the relevance of retrieved content and the quality of generated responses in the RAG framework by jointly leveraging the structure of knowledge graphs and the semantics of the query. Specifically, entities in a knowledge graph often have multiple adjacent paths, each representing different potential semantic associations. We begin by extracting key entities from the input query and use them as anchors to perform multi-path semantic expansion within the knowledge graph, thereby retrieving information that is more helpful for answering the original query. Considering the heterogeneity of semantic relations in different paths—such as one-hop, multi-hop, or importance-based paths—we design a multi-path subgraph construction mechanism to comprehensively capture relevant semantic evidence. However, many paths or triples may lack precise semantic alignment with the query, potentially introducing

irrelevant or noisy information into the retrieval process. To address this, we introduce a query-aware attention reward model, which performs fine-grained scoring of triples across different paths based on their semantic relevance to the query. Compared to traditional vector similarity-based methods, our attention model exhibits stronger query-awareness, thereby reducing noise introduced by semantic mismatches. We then select the most query-relevant subgraph as the backbone and fuse high-relevance triples from other paths to form a compact, semantically aligned fused subgraph. This fused subgraph is subsequently used to expand the query, further improving the relevance of document retrieval and ultimately enhancing the response quality of LLMs. We evaluate QMKGF on five benchmark datasets: SQuAD, IIRC, Culture, HotpotQA, and Musique. On the HotpotQA dataset, our method achieves a ROUGE-1 score of 64.98%, surpassing the BGE-Rerank approach by 9.72 percentage points (from 55.26% to 64.98%). Experimental results demonstrate that QMKGF significantly outperforms existing methods and shows clear advantages in improving answer accuracy.

In summary, the contributions of our paper are as follows:

1. A multi-path KG subgraph construction method is developed, incorporating one-hop relations, multi-hop relations, and importance-based relations to capture diverse and salient knowledge paths.
2. A query-aware attention reward model is introduced, enabling fine-grained scoring of subgraph triples based on their semantic relevance to the input query.
3. A subgraph fusion strategy is proposed, which selects the highest-scoring subgraph and integrates semantically relevant triples from other subgraphs, resulting in a more informative and query-aligned KG subgraph.
4. A KG-based query expansion approach is presented, enriching the original query with entities, relations, and triples from the final subgraph to improve the quality and factual consistency of RAG outputs.

2. Related Work

In recent years, Large Language Models (LLMs) have shown impressive performance in both comprehension and text generation, especially within the context of open-domain QA tasks [16, 25]. They have been widely applied in areas such as QA systems [29, 11, 32], code generation [17, 15, 1]. However, LLMs' heavy reliance on static knowledge embedded in their parameters. When applied to scenarios characterized by rapidly evolving information and high factual requirements, LLMs are prone to producing outdated or factually inaccurate responses. This limitation frequently leads to hallucinations, thereby undermining the overall quality and trustworthiness of the generated outputs [39, 38]. To mitigate the aforementioned limitations, Retrieval Augmented Generation (RAG) is a technique that combines information retrieval with text generation, aiming to improve the accuracy of generated content and reduce hallucinations

by incorporating support from external documents [9]. The typical workflow involves two stages: first, a retrieval model is used to fetch relevant information from external knowledge sources based on the user's query; then, this retrieved content is input into LLMs to integrate the information and generate a response.

RAG for LLMs

RAG has emerged as a mainstream approach for mitigating hallucination issues in LLMs due to its ability to retrieve relevant factual content. By introducing external factual knowledge during the inference process, RAG provides contextual support for LLMs, thereby enhancing the factual consistency and reliability of the generated content [3, 18, 34]. Existing LLM-based RAG methods can be broadly categorized into two paradigms: document retrieval-based RAG and graph retrieval-based RAG.

For document retrieval-based RAG, embedding serves as a critical component. It determines the ability to vectorize text, bringing semantically similar content closer in the vector space. The common method is to divide long text into several chunks, which weakens the semantic integrity and context coherence. To address this, Luo et al. [13] proposed Landmark Embedding, which employs a sliding-window embedding training method that preserves contextual consistency and enhances the model's ability to represent extended text. Additionally, to enrich the semantic content of user queries, Wang et al. [27] introduced a method that uses a few-shot prompting strategy to guide LLMs in generating pseudo-documents. These are concatenated with the original queries to form augmented queries, allowing the retrieval process to incorporate both the foundational knowledge of the LLMs and the content of external documents. Shi et al. [23] proposed Genground, a model that first uses the query as input to an LLM to generate a base answer, and then refines potential errors in this answer using retrieved content.

For graph retrieval-based RAG methods, graphs are typically integrated with the RAG framework to enhance the model's generalization capability [12]. Li et al. [5] proposed GraphRAG, which constructs semantic graphs to enable query-focused summarization, thereby improving the quality of responses to user queries. Zhu et al. [41] proposed KG2RAG, a framework that integrates structured knowledge from knowledge graphs to reinforce the semantic connections between document chunks, thereby improving both retrieval accuracy and generation quality. In the cultural tourism domain, Wei et al. [28] proposed QCG-rerank, which constructs a query-document graph structure and integrates query expansion with reranking mechanisms to enhance the relevance and expressiveness of retrieved passages.

Existing RAG methods are constrained by similarity-based matching between queries and chunks, often overlooking entities in the query and their deep semantic associations across chunks. To address this issue, we propose QMKGF, an approach that automatically constructs KGs. For the entities in the query, we first employ a multi-path KG subgraph construction strategy based on one-hop relations, multi-hop

relations, and importance-based relations. This approach enriches the relevance of entity information within the query, thereby enhancing the relevance of retrieved content during the recall process and ultimately improving the question-answering capabilities of LLMs.

3. Methodology

In this section, we provide a detailed introduction to QMKGF, whose overall architecture is illustrated in Fig. 2. The model consists of four main components: entity mapping, subgraph construction, subgraph fusion, and subgraph utilization. 1) For entity mapping, we first design prompts to leverage LLMs for KG extraction and construct an entity vector database. Then, potential entities are extracted from the query and mapped to real entities in the KG by matching them against the entity vector database (Section 3.1). 2) For subgraph construction, we propose a multi-path subgraph generation strategy that incorporates one-hop relations, multi-hop relations, and importance-based relations (Section 3.2). 3) For subgraph fusion, a query-aware attention reward model is used to score the three generated subgraphs. We then integrate valid triples from the lower-scoring subgraphs into the higher-scoring one to produce the final subgraph (Section 3.3). 4) In the subgraph utilization module, entities, relations, and triples from the final subgraph are concatenated with the original query for retrieval over a vector database. Relevant document chunks are retrieved and re-ranked using a reranking module, after which the top-ranked chunks are fed into the LLMs to generate the final response (Section 3.4).

3.1. Entity Mapping

In the context of processing a large collection of unstructured documents, denoted as $D = \{d_1, d_2, \dots, d_n\}$, we first leverage LLMs to extract entities and relations from the unstructured texts and construct a KG accordingly. The prompt design details for this construction are illustrated in Fig. 1.

Subsequently, we refined the embedding representations to improve their capacity for semantic comprehension and entity correspondence. The training dataset was constructed in the format $\{"query": str, "pos": List[str], "neg": List[str]\}$, where "pos" refers to the relevant answers, and "neg" includes irrelevant samples randomly drawn from other documents. This corpus was employed to adapt the embedding model to the target domain through fine-tuning. The optimized loss function is as follows:

$$\mathcal{L}_i = -\log \frac{\exp(\text{sim}(q_i, p_i)/m)}{\sum_j \exp(\text{sim}(q_i, d_j)/m)} \quad (1)$$

where, q_i denotes the embedding vector of the i -th query, p_i denotes positive document of the i -th query, d_j denotes all passages, including both positive and negative examples. m denotes the temperature scaling parameter. $\text{sim}()$ denotes the similarity calculation function.

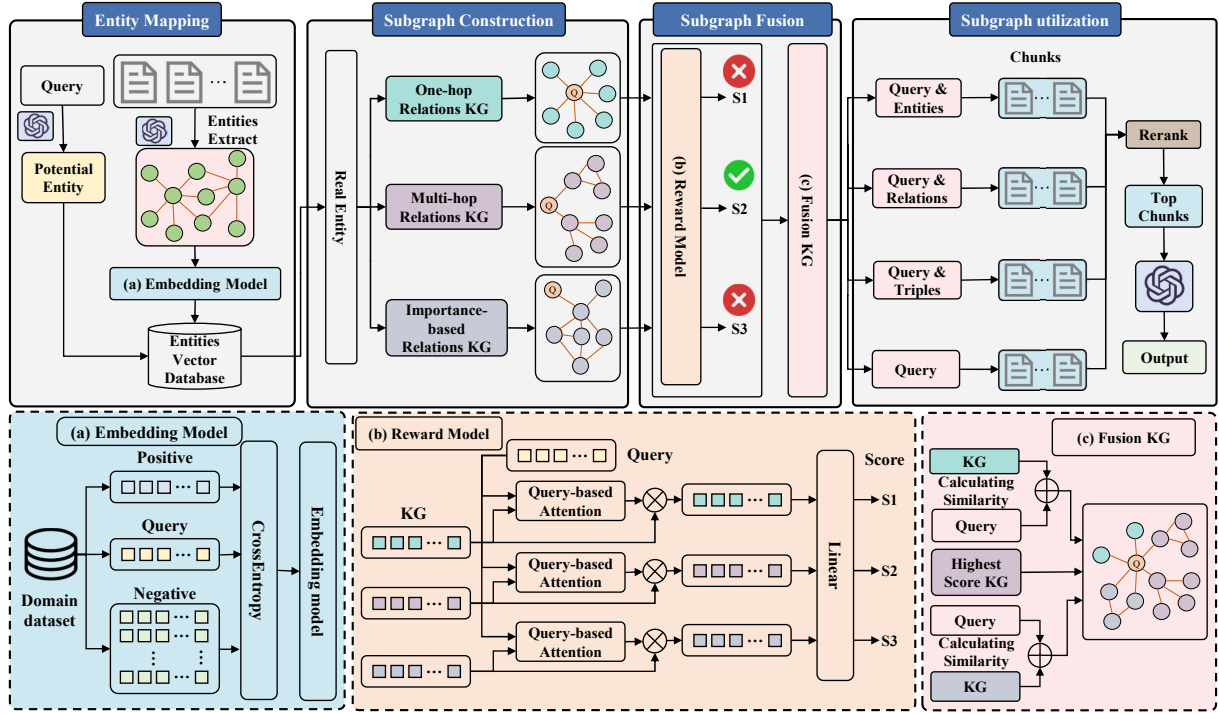


Figure 2: Framework of the proposed QMKGF.

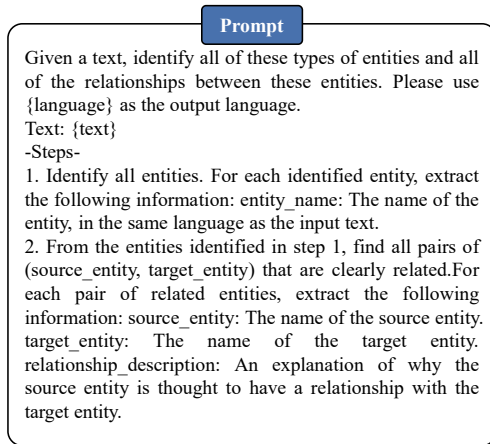


Figure 3: Prompt templates for entities and relations extraction.

Subsequently, we utilize the fine-tuned embedding model to vectorize the entities in the KG, resulting in entities vector database ent_VB to support entity matching in the subsequent retrieval process. Through the above steps, we ultimately obtain a complete KG $G=(V,E)$ and entities vector database ent_VB .

Given a query q , the entity mapping module first employs predefined prompt templates to extract potential entities e_p from the original query. Then, we calculate the similarity between e_p and each entity $e_i \in ent_VB$ to identify the most semantically similar entity. Through this process, the module

maps e_p to the most relevant existing entity $e_i \in ent_VB$ within the KG.

$$\text{sim}(e_p, e_i) = \frac{e_p \cdot e_i}{\|e_p\| \|e_i\|}, \quad \forall e_i \in ent_VB \quad (2)$$

3.2. Subgraph Construction

To fully leverage the contextual information associated with each entity, we design a multi-path subgraph construction algorithm. This algorithm integrates one-hop relations, multi-hop relations, and community-level entities based on personalized PageRank. The overall process is illustrated in Fig. 1.

One-hop relations: Given the identified truth entity e_t , we first construct its one-hop neighbor set $\mathcal{N}_1(e_t)$, defined as $\{v \in V \mid (e_t, v) \in E\}$. Within this set, we compute the semantic similarity between e_t and each neighboring entity $e_i \in \mathcal{N}_1(e_t)$, and select the top- K most relevant nodes. The resulting one-hop subgraph is defined as:

$$\text{Subg}_{\text{onehop}} = \text{top}_K(\text{sim}(e_t, e_i)), \quad e_i \in \mathcal{N}_1(e_t) \quad (3)$$

where $\text{sim}(\cdot)$ denotes the similarity computation function, and $\text{top}(\cdot)$ refers to selecting the top- K entities from a set. We select K one-hop neighbors of the target entity to enrich its representation by leveraging the contextual information provided by its immediate neighbors.

Multi-hop relations: We select the two most relevant entities $e_o \in \mathcal{N}_1(e_t)$ from the one-hop neighbors to perform multi-hop entity expansion. Specifically, we construct

a multi-hop entity set $NM = \{ v \in V \mid (e_o, v) \in E \}$, and then select the top- K most relevant nodes from this set to obtain the multi-hop subgraph, denoted as $Subkg_{mulhop}$.

Importance-based relations: In the importance scoring module, we employ the PageRank algorithm to quantify the relative importance of different components within the community. We first construct a personalized vector p in accordance with Eq. (4).

$$p(v) = \begin{cases} 1, & \text{if } v = e_t \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

By iteratively applying Eq. (5), we compute the PageRank score $S(v)$ for each node v :

$$S^{(t+1)}(v) = (1-d) \cdot p(v) + d \sum_{u \in In(v)} \frac{w_{uv}}{\sum_{k \in Out(u)} w_{uk}} S^{(t)}(u) \quad (5)$$

where d is the damping factor (commonly set to 0.85), $In(v)$ denotes the set of all nodes pointing to v , $Out(u)$ denotes the set of all nodes pointed to by node u , and w_{uv} , w_{uk} denote the weights of the respective edges.

We sort the nodes in descending order according to their PageRank scores $S(v)$, and select the top k nodes to obtain the subgraph $Subkg_{PR}$.

3.3. Subgraph Fusion

In the subgraph fusion module, we primarily accomplish two tasks: scoring the subgraphs using a query-aware attention reward model, and integrating the other subgraphs into the highest-scoring subgraph.

First, based on the three KG subgraphs constructed using the multi-path algorithm, we employ LLMs to score the dataset in terms of KG richness, question relevance, and connectivity. The scores assigned by the LLMs are used as training data to train query-aware attention reward model (RM).

To incorporate the semantic information of queries into the scoring process of the RM for KGs, we design a query-aware attention mechanism. We model the attention mechanism between the query vector and the representation of the KG subgraph, thereby emphasizing information regions that are highly relevant to the query semantics and ultimately enhancing the discriminative capability of the reward model.

Specifically, we first obtain the semantic representations of the query input $q \in \mathbb{R}^d$ and the KG subgraph $KGS \in \mathbb{R}^d$. Then, the query and subgraphs representations are linearly projected into the attention space to generate the Query, Key, and Value vectors used in the attention mechanism:

$$Q = q \cdot W_Q, K = KGS_i \cdot W_K, V = KGS_i \cdot W_V \quad (6)$$

where, $W_Q \in \mathbb{R}^{d \times d}$, $W_K \in \mathbb{R}^{d \times d}$, and $W_V \in \mathbb{R}^{d \times d}$ denote trainable weight matrices.

We employ a query-aware attention mechanism to compute the query-aware response representation of the subgraph:

$$\text{head} = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \quad (7)$$

$$\text{Attention}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \quad (8)$$

where W^O denotes trainable weight matrices, $\text{Concat}()$ denotes concatenation function.

This module design ensures that the reward computation not only depends on the inherent representation of the KG subgraph but also reflects the degree of semantic alignment with the query, thereby providing improved contextual adaptability and semantic guidance.

Subsequently, we apply RM to score the three KG subgraphs and select the highest-scoring subgraph. The selection process is described by Eq. (9):

$$KGS_{max} = \max \left(RM(KGS_{onehop}), RM(KGS_{mulhop}), RM(KGS_{PR}) \right) \quad (9)$$

where $\max()$ denotes selecting the highest score from the RM's outputs and returning the subgraph with the top score.

For subgraphs generated from different sources, although some subgraphs may have lower scores, they can still contain important triples highly relevant to the query entity. To fully exploit this potential valuable information, we design a similarity-based subgraph fusion method. We compute the similarity between the highest-scoring subgraph and the query, setting this similarity as a threshold r . The specific calculation formula is as follows:

$$r = \cos(KGS_{max}, q) \quad (10)$$

For each lower-scoring subgraph $LS_i=(h,r,t)$, we calculate the similarity between each triple and the query as follows:

$$\text{sim}((h, r, t), q) = \cos((h, r, t), q) \quad (11)$$

If a triple satisfies Eq. (12), it is considered highly relevant to the query and should be retained. All such triples are set as $T_{selected}$.

$$\text{sim}((h, r, t), q) \geq \tau \quad (12)$$

Finally, the fused subgraph G_{fusion} is obtained by Eq. (13):

$$G_{fusion} = Subkg_{max} \cup T_{selected} \quad (13)$$

3.4. Subgraph Utilization

In the process of query-aware retrieval, model performance is often limited by the narrow scope of information coverage and insufficient semantic relevance. To address this issue, we propose a query expansion mechanism enhanced with a KG. It works together with G_{fusion} to improve retrieval relevance and support the performance of RAG. Specifically, we extract an entity set e_{fusion} , a relation set

r_{fusion} , and a triple set T_{fusion} from G_{fusion} . The query is concatenated with each entity, relation and triple to form distinct retrieval items, which are then used to perform similarity searches in the vector database. The resulting document sets are aggregated into the final collection Doc .

To further enhance the alignment between retrieved passages and the query task, we employ the bge-rerank re-ranking model to score all candidate documents in Doc based on their relevance. The top- k documents are then selected according to the ranking results. This process is formally defined in Eq. (14).

$$C_{rank} = \text{top}(\text{rerank}(Doc)) \quad (14)$$

Subsequently, the filtered set C_{rank} is used as contextual input to the LLMs, which summarize and synthesize the information from the selected passages to generate the final answer.

$$\text{Output} = \text{LLM}(C_{rank}) \quad (15)$$

3.5. QMKGF algorithmic

Through the above process, we obtain the response generated by LLMs based on chunks filtered and expanded using the QMKGF algorithm, which incorporates semantic information from KG subgraphs into the original query. The overall algorithm workflow is as shown in Table 1.

4. Experiment Setup

In this section, we introduce the datasets, baseline models, experimental settings, and metrics.

4.1. Datasets

We used Cultour, IIRC, HotpotQA, SQuAD and MuSiQue datasets to evaluate the models.

Cultour [29] is a tourism-oriented dataset comprising 12,000 QA pairs, sourced from both manually curated travel-related materials and content automatically generated by large language models. It is employed to assess model performance in domain-specific tourism scenarios.

IIRC [6] is built from English Wikipedia and includes over 13,000 questions. Each question offers only partial context, requiring models to retrieve and synthesize information scattered across multiple related passages. This dataset serves to test comprehension under conditions of incomplete or distributed information.

HotpotQA [36] consists of around 113K questions, where correct answers demand the integration of facts drawn from two distinct Wikipedia articles. It is used to benchmark a model's capability in cross-document reasoning and complex inference.

SQuAD [19] is a benchmark reading comprehension dataset where answers are typically embedded within a designated paragraph of a Wikipedia entry. It evaluates how well a model can extract and comprehend localized textual information.

Table 1

The overall QMKGF algorithm.

Algorithm QMKGF

Input Query q , Documents chunks, LLMs

Output LLMs output

1. Begin
2. Construct KG using LLMs
3. Fine-tune embedding on domain datasets
4. $document_VB \leftarrow \text{document}$, $ent_VB \leftarrow \text{entity}$
5. For q in query sets:
6. Extract potential entities e_p from the query
7. $e_t \leftarrow \text{mapping of } e_p \text{ in KG}$
8. Construct multiple subgraphs based on one-hop relations $Subkg_{onehop}$, multi-hop relations $Subkg_{mulhop}$, and importance-based relations $Subkg_{PR}$ of e_t
9. Design a reward model function based on query-aware attention mechanism
10. $Q \leftarrow q \cdot W_Q$, $K \leftarrow KGS_i \cdot W_K$, $V \leftarrow KGS_i \cdot W_V$
11. $\text{head} \leftarrow \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$
12. $\text{Attention}(Q, K, V) \leftarrow \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$
13. $\text{score} \leftarrow \text{linear}(\text{Attention})$
14. $KGS_{max} \leftarrow \max(RM(KGS_{onehop}, KGS_{mulhop}, KGS_{PR}))$
15. $r \leftarrow \cos(KGS_{max}, q)$
16. $T_{selected} \leftarrow \{(h, r, t)\} \text{ where } \text{sim}((h, r, t), q) \geq \tau$
17. $G_{fusion} \leftarrow Subkg_{max} \cup T_{selected}$
18. Semantically expand the query q based on G_{fusion}
19. Retrieve relevant documents Doc
20. $C_{rank} \leftarrow \text{top}(\text{rerank}(Doc))$
21. Input the top k chunks C_{rank} into the LLMs to generate output
22. End begin

MuSiQue [26] focuses on multi-hop reasoning, featuring questions that generally require 2–4 logical steps to reach a conclusion. It is utilized to assess the model's capacity for handling multi-hop inference tasks across multiple supporting contexts.

4.2. Baseline

In this section, we introduce the baseline models in detail.

LLMs-only: The query is directly input into the LLMs without retrieving any relevant documents from an external database. The final answer is generated solely by the LLMs.

W-RAG: Based on the input query, relevant information is retrieved from the database and appended to the query as supplementary context, which is subsequently processed by the LLMs to produce the final answer.

BM25 [21] is a statistics-based information retrieval method that scores the relevance between documents and queries by combining the term frequency of query words within a document and the inverse document frequency of those terms across the entire corpus.

BGE-rerank [33] first performs an initial retrieval based on semantic similarity. Subsequently, the BGE-rerank model is applied to re-rank the retrieved results, thereby improving the relevance of the documents input into LLMs.

BCE-rerank [14], developed by NetEase Youdao, is a bilingual and cross-lingual semantic embedding model specialized in enhancing semantic search accuracy and refining

Table 2

The overall experimental results of QMKGF and other baselines on SQuAD, IIRC, Cultour datasets. The best results are in bold.

Models	SQuAD				IIRC				Cultour			
	R-1	R-L	B-1	Met.	R-1	R-L	B-1	Met.	R-1	R-L	B-1	Met.
LLMs-only	8.17	7.96	20.12	14.12	5.09	5.06	14.16	11.49	26.51	18.72	35.95	20.09
W-RAG	40.58	40.46	46.63	38.70	45.74	45.68	52.61	37.29	60.13	55.68	49.78	44.57
BM25	43.63	43.47	48.81	41.59	47.27	47.19	53.04	38.38	61.17	56.44	50.74	45.34
BGE-rerank	43.47	43.36	48.51	41.65	46.41	46.38	53.43	37.42	61.23	56.56	50.66	45.65
BCE-rerank	43.25	43.13	48.73	41.67	45.37	45.43	52.29	37.16	61.52	56.93	50.97	45.73
LLMs-KG	8.76	8.53	19.05	13.89	8.09	8.02	15.19	12.20	30.36	21.67	43.48	23.14
LLMs-KG-rerank	43.86	43.74	49.44	41.10	46.41	46.26	54.42	37.92	60.44	55.82	50.35	44.64
QCG-rerank	45.32	45.12	50.95	42.56	47.11	47.05	55.13	38.62	61.89	57.12	51.41	45.49
QMKGF(ours)	51.44	51.26	56.26	47.75	50.37	50.33	57.53	40.38	62.61	58.55	52.12	45.96

the ranking order of search results according to semantic relevance.

LLMs-KG performs semantic retrieval on an automatically generated KG based on entities extracted from the query, and incorporates the retrieved information into the input of the LLMs to enhance the quality of the generated output.

LLMs-KG-rerank extracts entities from the query and retrieves related chunks from an automatically constructed KG. The retrieved chunks are then re-ranked using the BGE-Rerank model, and the top-ranked chunks is input into the LLMs for answer generation.

KG2RAG [41] employs a KG-guided chunk expansion process and a KG-based chunk organization process to deliver relevant and important knowledge in well-organized paragraphs.

QCG-rerank [28] introduces query expansion and chunk graph re-ranking mechanisms. It enriches the semantic representation of the query by extracting and duplicating key information, then constructs a semantic similarity-based chunk graph for iterative re-ranking. The final ranked results are subsequently input into LLMs.

4.3. Experiment settings

To ensure experimental fairness, we selected qwen2.5-7B-Instruct as the base LLM and set the temperature to 0.0. During embedding fine-tuning, the number of training epochs was uniformly set to 10, with a maximum query length of 64 tokens and a maximum passage length of 256 tokens. The learning rate was fixed at $1e-5$. For the re-ranking stage, we employed the bge-rerank model, utilizing the BGE Large embedding for vector representations. For the reward model (RM), we used ernie-3.0-base with a learning rate of $1e-5$, a maximum input length of 512 tokens, and trained for 10 epochs. Entity and relation extraction as well as KG construction for Chinese datasets were performed using the qwen2.5 model, while for English datasets, Llama 3 was used. For the IIRC dataset, 4,906 samples were used for training and 593 for testing. The Cultour dataset was split into training and testing sets with an 8:2 ratio. For HotpotQA, SQuAD, and MuSiQue datasets, the same 8:2 split was applied. If the test set contains more than 1,000

samples, a random subset of 1,000 samples is selected for evaluation. All experiments were conducted on an NVIDIA RTX A6000 GPU.

4.4. Metrics

To evaluate model performance, we adopt three metrics: ROUGE, BLEU, and METEOR. ROUGE [2] captures lexical overlap between generated outputs and human-annotated references, emphasizing recall. BLEU [20] focuses on precision by computing n-gram co-occurrence between candidate and reference sequences. METEOR [22] provides a more fine-grained evaluation by accounting for not only exact word matches but also stemming, synonymy, and word order, thus offering a more comprehensive assessment of generation quality compared to BLEU.

5. Results

In this section, we evaluate the proposed QMKGF model through a comparative analysis against a number of representative baselines. Subsequently, ablation studies are conducted to investigate the impact of key components and design choices within our framework.

5.1. Main results

We evaluated QMKGF on the SQuAD, IIRC, Cultour, HotpotQA, and MuSiQue datasets, and the detailed results are presented in Table 2 and Table 3.

As shown in Tables 2 and 3, our proposed QMKGF framework consistently outperforms all baseline models across five benchmark datasets: SQuAD, IIRC, Cultour, HotpotQA, and MuSiQue, with the best performance highlighted in bold. In particular, QMKGF achieves the highest scores in nearly all evaluation metrics, such as ROUGE-1, ROUGE-L, BLEU-1, and METEOR. The most notable gain is on the HotpotQA dataset, where QMKGF surpasses bge-rerank by 9.72. QMKGF attains a remarkable 64.98 ROUGE-1 and 64.95 ROUGE-L, outperforming the previous best (KG2RAG) by over 6.0 points. Similar results were achieved on two other multi-hop question answering datasets, IIRC and MuSiQue, demonstrating its superiority in multi-hop reasoning tasks. QMKGF achieves significant

Table 3

The overall experimental results of QMKGF and other baselines on HotpotQA and MuSiQue datasets. The best results are in bold.

Models	HotpotQA				MuSiQue			
	R-1	R-L	B-1	Met.	R-1	R-L	B-1	Met.
LLMs-only	11.86	11.73	25.38	17.91	2.49	2.43	13.20	10.75
W-RAG	54.07	54.01	60.22	47.20	36.23	36.16	42.82	34.89
BM25	54.86	54.73	60.76	50.06	38.53	38.20	43.64	36.36
BGE-rerank	55.26	55.15	61.30	49.91	39.51	39.43	45.50	37.05
BCE-rerank	54.97	54.87	60.45	49.43	39.13	39.11	45.37	36.83
LLMs-KG	12.59	12.52	25.11	17.77	2.89	2.82	11.70	9.99
LLMs-KG-rerank	54.16	54.06	60.04	48.50	38.82	38.73	44.64	36.08
QCG-rerank	57.44	57.24	62.69	51.61	41.17	40.98	46.53	38.31
KG2RAG	58.90	58.83	66.65	53.32	37.41	37.09	48.25	35.71
QMKGF(ours)	64.98	64.95	68.42	57.74	47.42	47.35	53.31	43.95

Table 4

Ablation experiments of QMKGF. The best results are in bold.

Models	HotpotQA				Cultour			
	R-1	R-L	B-1	Met.	R-1	R-L	B-1	Met.
Bge-rerank	55.26	55.15	61.30	49.91	61.23	56.56	50.66	45.65
w/o-attention	61.17	61.13	65.79	54.57	61.76	58.03	49.86	45.10
w/o-fintune	61.20	61.15	65.83	54.90	61.77	58.06	49.69	45.13
QMKGF	64.98	64.95	68.42	57.74	62.61	58.55	52.12	45.96

improvements on the Chinese question answering dataset Cultour. Specifically, it outperforms the rerank-only baseline bge-rerank by 1.39 points on the ROUGE-1 metric. On the English extraction dataset SQuAD, QMKGF is 7.97 points higher than the BGE-rerank model in the ROUGE-1 metric. This demonstrates the model's superiority in single-hop question answering tasks.

Moreover, across all datasets, QMKGF consistently outperforms both the LLMs+KG and LLMs+KG+rerank. These results indicate that QMKGF effectively enhances the quantity of relevant chunks by leveraging query-aware attention RM to filter and integrate multiple subgraph paths. The reranked results are then passed to LLMs, leading to improvements in both the relevance and reliability of the generated responses.

5.2. Ablation experiment

To assess the effectiveness of our proposed approach, experiments were conducted on the Chinese Cultour dataset and the English HotpotQA dataset. The main tests focused on the impact of the following components on the model's performance: without query-aware attention (w/o-attention), and without fine-tuned embedding (w/o-fintune). The detailed results are shown in Table 4.

As shown in Table 4, each component contributes positively to the final performance. QMKGF achieves the best performance on R-1, R-L, and B-1, demonstrating the effectiveness of the proposed multi-path KG subgraph construction strategy. In the results on the English dataset HotpotQA, we observe that all components of our framework lead to

a substantial improvement over the BGE-rerank baseline. While similar gains are also present on the Chinese dataset, the improvement is less pronounced compared to IIRC. Additionally, without personalized PageRank algorithm has a more significant negative impact on English datasets than on the Chinese dataset. One possible reason is that we adopt the Qwen2.5 model as the backbone, which has been extensively pretrained on Chinese corpora, leading to stronger prior knowledge in the Chinese domain. In contrast, its relative lack of exposure to English knowledge makes the inclusion of high-quality triples from diverse subgraphs more influential for final performance on English tasks.

The performance impact of removing query-aware attention and removing the fine-tuned embedding shows a similar trend. These results suggest that QMKGF effectively integrates one-hop relations, multi-hop relations, and the importance-based relations algorithm to construct and fuse multiple semantic paths. In the process of scoring the three subgraphs using the query-aware attention reward model, the influence of the query on each subgraph is taken into account. Thereby improving the quality of the retrieved documents and enhancing the performance of LLMs in response generation.

5.3. Impact of fine-tuning embeddings

To evaluate the impact of fine-tuned embedding models on retrieval performance, we conducted experiments using metrics such as MRR@1, MRR@10, Recall@10, and nDCG@10. As illustrated in Fig. 4, the fine-tuned models show consistent and substantial improvements across all

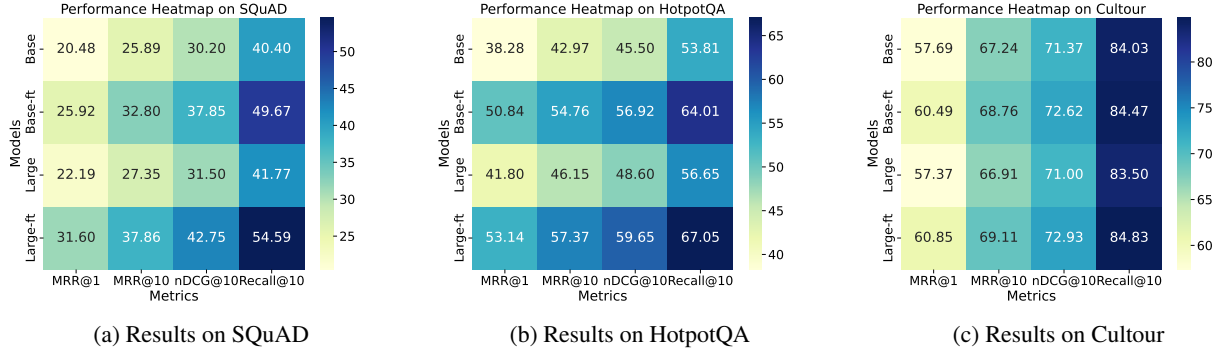


Figure 4: Impact of fine-tuning embeddings (combined view).

Table 5

Results on Precision, Recall, and F1-score. The best results are in bold.

Models	Cultour			
	P	R	F1	hit@10
W-RAG	63.44	53.61	58.12	97.70
BGE-rerank	62.03	53.43	57.41	98.00
QMKGf(ours)	71.60	51.53	59.24	99.40

metrics on both English datasets (SQuAD, HotpotQA) and the Chinese dataset (Cultour). For example, on the SQuAD dataset, fine-tuning the bge_base model results in a 9.27-point increase in Recall@10, while the bge_large model achieves an even greater improvement of 12.82 points. This indicates that fine-tuning the bge_large model has a greater impact on the final performance. Moreover, we observe that the performance improvement from embedding fine-tuning is more pronounced on the English datasets compared to the Chinese dataset. This is likely because the Chinese dataset already achieved relatively strong performance before fine-tuning, leaving limited room for further improvement, which in turn makes the fine-tuning gains less noticeable. These results demonstrate that domain-specific embedding fine-tuning can effectively enhance the accuracy of relevant document retrieval.

5.4. Assessment of Information Retrieval Ability

In addition to evaluating text quality using ROUGE, BLEU, and METEOR, we also incorporate Hit@10, as well as Precision, Recall, and F1-score, to provide a complementary assessment of QMKGf's generated results from the perspectives of answer matching and information retrieval. These metrics offer a more comprehensive evaluation of the relevance and accuracy of the generated content.

As shown in Table 5, QMKGf outperforms both the W-RAG and BGE-rerank models. In the hit@10 metric, QMKGf achieved 99.4%, reaching the highest performance, indicating that our model has a significant advantage in the accuracy of answer retrieval and matching. In terms of Precision (P) and F1-score(F1), our model achieved the best

performance, demonstrating excellent accuracy in generating results. The results indicate that our multi-path knowledge graph approach effectively expands the semantic neighborhood of query entities, broadens the recall scope, and enhances query relevance, leading to notable improvements in retrieval and generation performance.

5.5. Impact of reward model

To further evaluate the impact of using different pre-trained models as reward models in QMKGf, we conducted comparative experiments with BERT [4] and ERNIE [24]. Both models possess strong semantic understanding capabilities and have been widely adopted in various NLP tasks. As shown in Table 6, employing ERNIE as the RM for subgraph selection leads to greater improvements in the accuracy of the LLM's outputs, demonstrating its superior effectiveness in guiding the retrieval process.

These improvements can be attributed to the semantic knowledge integration mechanisms of ERNIE, which leverages entity-level masking and prior knowledge during pretraining. Unlike BERT, which is limited to token-level representations, ERNIE incorporates structured knowledge from knowledge graphs and factual information during its representation learning process. This enables the reward model to better evaluate the semantic relevance between the candidate triples and the query, thereby guiding the selection of more informative subgraphs and enhancing the overall output quality.

5.6. Effect of Subgraph Node Count

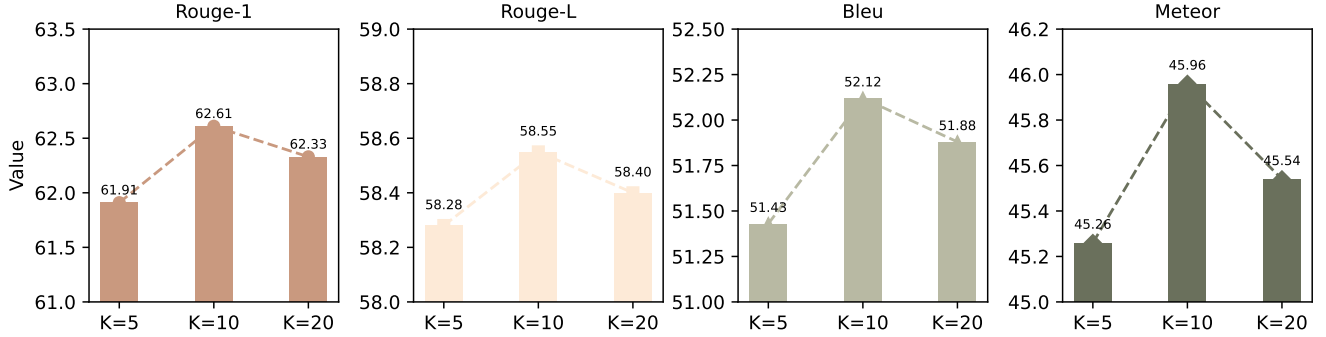
To evaluate the impact of different subgraph sizes on the generation performance of QMKGf, we conducted comparative experiments on Cultour datasets with three configurations, setting the number of subgraph nodes K to 5, 10, and 20, respectively.

As shown in Fig. 5, the model achieves optimal performance when $K=10$, with the highest scores across all evaluation metrics, including ROUGE-1 (62.61) and BLEU (52.12). This indicates that the configuration with $K=10$ strikes a favorable balance between semantic relevance and information redundancy, thereby effectively enhancing generation quality. When the number of nodes is smaller ($K=5$),

Table 6

The impact of reward model. The best results are in bold.

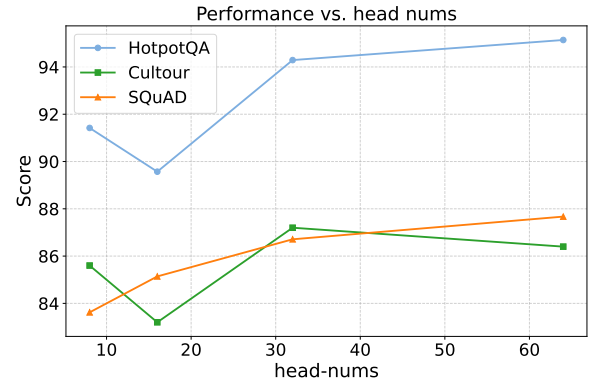
Models	Cultour				SQuAD			
	R-1	R-L	B-1	Met.	R-1	R-L	B-1	Met.
ERNIE	62.61	58.55	52.12	45.96	51.44	51.26	56.26	47.75
BERT	62.06	58.40	49.97	45.30	51.05	50.86	55.89	47.45

**Figure 5:** Effect of subgraph node count.

the performance on ROUGE metrics remains relatively close, but BLEU and METEOR scores are notably lower. This suggests insufficient semantic coverage in the subgraph, which results in less informative and under-supported generated content. On the other hand, increasing the node count to $K=20$ leads to a slight decline in performance, particularly in BLEU and METEOR scores. This decline may be attributed to the introduction of irrelevant or noisy entities and relations, which interferes with semantic alignment and coherent context construction during generation.

5.7. The impact of head nums on RM performance

Fig. 6 illustrates the impact of the number of attention heads in the multi-head attention mechanism on model performance, showing notable differences across the HotpotQA, Cultour, and SQuAD datasets. Experimental results demonstrate that increasing the number of attention heads significantly improves performance on HotpotQA, with the highest scores observed at 32 and 64 heads. This suggests that for complex multi-hop reasoning tasks, introducing more parallel attention paths facilitates the capture of diverse information sources and the modeling of long-range dependencies, thereby enhancing the model's ability to identify semantically relevant knowledge paths. In contrast, the performance improvements on Cultour and SQuAD are relatively modest. This may be attributed to the relatively lower reasoning complexity and simpler information structures of these tasks, where an excessive number of attention heads could introduce redundant computations or semantic noise, thus weakening the model's focus on core information.

**Figure 6:** The impact of head nums on RM performance.

These results indicate that increasing the number of attention heads is highly beneficial for complex tasks such as multi-hop question answering, as it enhances the model's capacity to capture and integrate knowledge paths. However, for tasks with more concentrated and structurally straightforward information, a moderate number of attention heads may already suffice for high-quality modeling, and further increases could lead to performance instability.

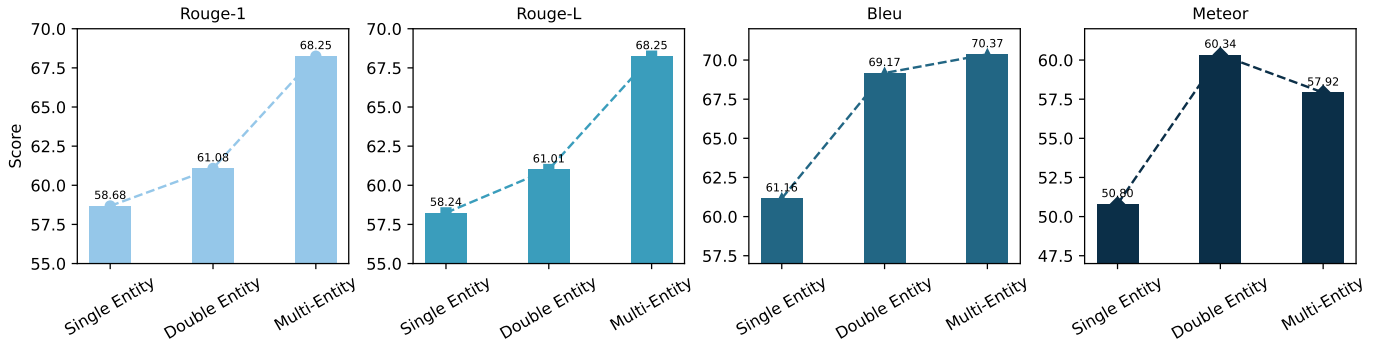
5.8. Impact of different subgraph fusion strategies

To further examine the impact of different subgraph fusion strategies on model performance, we conduct experiments comparing three representative approaches—All

Table 7

Impact of different subgraph fusion strategies. The best results are in bold.

Models	HotpotQA				MuSiQue			
	R-1	R-L	B-1	Met.	R-1	R-L	B-1	Met.
All fusion	64.07	64.05	67.62	57.01	46.36	46.28	52.42	43.24
Top-5 triples fusion	63.54	63.51	67.26	56.84	43.33	44.24	50.93	41.67
RM fusion	64.98	64.95	68.42	57.74	47.42	47.35	53.31	43.95

**Figure 7:** Effect of subgraph node count.

fusion, Top-5 triples fusion, and RM fusion—on the HotpotQA and MuSiQue datasets. All Fusion: All triples from each subgraph are retained without any filtering. Top-5 Triples Fusion: Only the top 5 triples from each subgraph are preserved to reduce noise. RM Fusion: A reward model is used to score and filter triples, retaining only those highly relevant to the query to enhance information quality and generation performance. The generation quality is evaluated using four standard metrics: ROUGE-1 (R-1), ROUGE-L (R-L), BLEU-1 (B-1), and METEOR (Met.). The experimental results are shown in Table 7.

Different subgraph fusion methods have a significant impact on the quality of the final generated results. Notably, the Top 5 triples fusion method performs well on the HotpotQA dataset but poorly on the MuSiQue dataset. The RM fusion method achieves the best results on both HotpotQA and MuSiQue, demonstrating that this strategy has a certain degree of generalizability. This is because All fusion may introduce excessive redundant or noisy information, which interferes with the model's focus on key content. Top-5 triples fusion reduces noise to some extent, but its fusion strategy is relatively coarse. It fails to effectively model semantic correlations, resulting in fewer relevant triples and thus limited document retrieval. In contrast, RM fusion models semantic relevance across multiple heterogeneous subgraphs more effectively. This helps retain the most valuable information, enables the retrieval of more relevant documents, and ultimately improves the quality of generated content.

5.9. Impact of entity quantity in queries

To further investigate the impact of entity quantity in queries on model performance, we divide the HotpotQA test set into three categories: single-entity queries, two-entity queries, and multi-entity queries (more than 2 entities). For each category, 200 samples are randomly selected for evaluation. The experimental results are shown in Fig. 7. The experimental results demonstrate a clear upward trend in the scores of the QMKGF model across three evaluation metrics—ROUGE-L, ROUGE-1, and BLEU—as the number of entities in the query increases, indicating enhanced generation capabilities. For instance, under the BLEU metric, the score improves significantly from 61.16% in the single-entity group to 70.37% in the multi-entity group.

This phenomenon reflects that, in the context of the multi-hop question answering task represented by HotpotQA, multi-entity queries provide more explicit semantic anchors, which facilitate richer semantic expansion when integrated with knowledge graphs. Consequently, the model is able to retrieve more query-relevant passages from the database during the retrieval stage, offering stronger contextual support for generation. In contrast, single-entity queries often exhibit semantic ambiguity, which may lead to overly broad retrieval scopes or insufficient contextual aggregation, thereby undermining the accuracy and consistency of the generated content.

6. Conclusion

In this paper, we proposed QMKGF, a novel framework aimed at enhancing the performance of large language models (LLMs). To capture semantically rich structures associated with the query, a multi-path subgraph construction strategy was developed. Subsequently, a query-aware attention reward model was introduced to score triples based on their semantic relevance to the query, guiding the selection of high-quality triples for knowledge graph fusion. Finally, the updated subgraph was utilized to expand the original query, thereby improving the relevance of retrieved documents. The proposed QMKGF was evaluated on multiple benchmark datasets, including SQuAD, IIRC, Culture, HotpotQA, and MuSiQue, and the experimental results demonstrate its effectiveness. Detailed ablation studies further validate the contributions of the domain-adaptation KG construction module, the query-aware attention reward model, and the subgraph fusion module. Additionally, we investigated the impact of fine-tuned embedding models on overall performance. Different pretrained models were evaluated as the backbone of the reward model. In addition, the impact of varying the number of heads in the query-aware attention mechanism on the final performance was also explored. Finally, we analyzed the influence of subgraph size in the knowledge graph on model performance.

Future work will explore the construction of more domain-specific knowledge graphs tailored to the unique characteristics of specialized data. Based on these KGs, the integration of causal relationship modeling and event chain

References

- [1] Bai, X., Huang, S., Wei, C., Wang, R., 2025. Collaboration between intelligent agents and large language models: A novel approach for enhancing code generation capability. *Expert Systems with Applications* 269, 126357. doi:<https://doi.org/10.1016/j.eswa.2024.126357>.
- [2] Barbella, M., Tortora, G., 2022. Rouge metric evaluation for text summarization techniques. Available at SSRN 4120317.
- [3] Cui, J., Li, Z., Yan, Y., Chen, B., Yuan, L., 2023. Chatlaw: Open-source legal large language model with integrated external knowledge bases. *arXiv preprint arXiv:2306.16092*.
- [4] Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. 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.
- [5] Edge, D., Trinh, H., Cheng, N., Bradley, J., Chao, A., Mody, A., Truitt, S., Metropolitansky, D., Ness, R.O., Larson, J., 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*.
- [6] Ferguson, J., Gardner, M., Hajishirzi, H., Khot, T., Dasigi, P., 2020. Iirc: A dataset of incomplete information reading comprehension questions. *arXiv preprint arXiv:2011.07127*.
- [7] Fu, Y., Liu, D., Zhang, B., Jiang, Z., Mei, H., Guan, J., 2025. Cue rag: Dynamic multi-output cue memory under h framework for retrieval-augmented generation. *Neurocomputing*, 130235.
- [8] Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Wang, M., Wang, H., 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- [9] Guu, K., Lee, K., Tung, Z., Pasupat, P., Chang, M., 2020. Retrieval augmented language model pre-training, in: *International conference on machine learning*, PMLR. pp. 3929–3938.
- [10] Jiang, Z., Xu, F.F., Gao, L., Sun, Z., Liu, Q., Dwivedi-Yu, J., Yang, Y., Callan, J., Neubig, G., 2023. Active retrieval augmented generation, in: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 7969–7992.
- [11] Li, X., Hou, X., Ravi, N., Huang, Z., Gan, Y., 2025. A two-stage proactive dialogue generator for efficient clinical information collection using large language model. *Expert Systems with Applications*, 127833.
- [12] Li, X., Wu, Z., Wu, J., Cui, H., Jia, J., Li, R.H., Wang, G., 2024. Graph learning in the era of llms: A survey from the perspective of data, models, and tasks. *arXiv preprint arXiv:2412.12456*.
- [13] Luo, K., Liu, Z., Xiao, S., Liu, K., 2024. Bge landmark embedding: A chunking-free embedding method for retrieval augmented long-context large language models. *arXiv preprint arXiv:2402.11573*.
- [14] NetEase Youdao, I., 2023. Bcembedding: Bilingual and crosslingual embedding for rag. <https://github.com/netease-youdao/BCEmbedding>.
- [15] Nie, L.Y., Gao, C., Zhong, Z., Lam, W., Liu, Y., Xu, Z., 2021. Coregen: Contextualized code representation learning for commit message generation. *Neurocomputing* 459, 97–107.
- [16] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al., 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35, 27730–27744.
- [17] Patsakis, C., Casino, F., Lykousas, N., 2024. Assessing llms in malicious code deobfuscation of real-world malware campaigns. *Expert Systems with Applications* 256, 124912. doi:<https://doi.org/10.1016/j.eswa.2024.124912>.
- [18] Peng, B., Galley, M., He, P., Cheng, H., Xie, Y., Hu, Y., Huang, Q., Liden, L., Yu, Z., Chen, W., et al., 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*.
- [19] Rajpurkar, P., 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- [20] Reiter, E., 2018. A structured review of the validity of bleu. *Computational Linguistics* 44, 393–401.
- [21] Robertson, S., Zaragoza, H., et al., 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 333–389.
- [22] Saadany, H., Orasan, C., 2021. Bleu, meteor, bertscore: Evaluation of metrics performance in assessing critical translation errors in sentiment-oriented text. *arXiv preprint arXiv:2109.14250*.
- [23] Shi, Z., Zhang, S., Sun, W., Gao, S., Ren, P., Chen, Z., Ren, Z., 2024. Generate-then-ground in retrieval-augmented generation for multi-hop question answering. *arXiv preprint arXiv:2406.14891*.
- [24] Sun, Y., Wang, S., Li, Y., Feng, S., Chen, X., Zhang, H., Tian, X., Zhu, D., Tian, H., Wu, H., 2019. Ernie: Enhanced representation through knowledge integration. *arXiv preprint arXiv:1904.09223*.
- [25] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al., 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- [26] Trivedi, H., Balasubramanian, N., Khot, T., Sabharwal, A., 2022. ♪ musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics* 10, 539–554.
- [27] Wang, L., Yang, N., Wei, F., 2023. Query2doc: Query expansion with large language models. *arXiv preprint arXiv:2303.07678*.
- [28] Wei, Q., Yang, M., Han, C., Wei, J., Zhang, M., Shi, F., Ning, H., 2024a. Qcg-rerank: Chunks graph rerank with query expansion in retrieval-augmented llms for tourism domain. *arXiv preprint arXiv:2411.08724*.
- [29] Wei, Q., Yang, M., Wang, J., Mao, W., Xu, J., Ning, H., 2024b. Tourllm: Enhancing llms with tourism knowledge. *arXiv preprint arXiv:2407.12791*.
- [30] Wu, J., Wu, Z., Li, R., Qin, H., Wang, G., 2024a. Effective bug detection in graph database engines: An llm-based approach. *arXiv preprint arXiv:2402.00292*.

- [31] Wu, J., Zhu, J., Qi, Y., 2024b. Medical graph rag: Towards safe medical large language model via graph retrieval-augmented generation. arXiv preprint arXiv:2408.04187 .
- [32] Xiao, Q., Li, R., Yang, J., Chen, Y., Jiang, S., Wang, D., 2024. Tpke-qa: A gapless few-shot extractive question answering approach via task-aware post-training and knowledge enhancement. *Expert Systems with Applications* 254, 124475.
- [33] Xiao, S., Liu, Z., Zhang, P., Muennighof, N., 2023. C-pack: Packaged resources to advance general chinese embedding. arXiv preprint arXiv:2309.07597 .
- [34] Xu, S., Pang, L., Shen, H., Cheng, X., Chua, T.s., 2023. Search-in-the-chain: Towards the accurate, credible and traceable content generation for complex knowledge-intensive tasks. arXiv preprint arXiv:2304.14732 .
- [35] Xu, T., Chen, L., Hu, Z., Li, B., 2025. Staf-llm: A scalable and task-adaptive fine-tuning framework for large language models in medical domain. *Expert Systems with Applications* 281, 127582. doi:<https://doi.org/10.1016/j.eswa.2025.127582>.
- [36] Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W.W., Salakhutdinov, R., Manning, C.D., 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. arXiv preprint arXiv:1809.09600 .
- [37] Yu, H., Gan, A., Zhang, K., Tong, S., Liu, Q., Liu, Z., 2024. Evaluation of retrieval-augmented generation: A survey, in: *CCF Conference on Big Data*, Springer. pp. 102–120.
- [38] Zhang, Y., Zheng, W., Huang, J., Xiao, G., 2025. Lgkgr: A knowledge graph reasoning model using llms augmented gnns. *Neurocomputing* 635, 129919.
- [39] Zhao, P., Zhang, H., Yu, Q., Wang, Z., Geng, Y., Fu, F., Yang, L., Zhang, W., Cui, B., 2024a. Retrieval-augmented generation for ai-generated content: A survey. arXiv preprint arXiv:2402.19473 .
- [40] Zhao, P., Zhang, H., Yu, Q., Wang, Z., Geng, Y., Fu, F., Yang, L., Zhang, W., Jiang, J., Cui, B., 2024b. Retrieval-augmented generation for ai-generated content: A survey. arXiv preprint arXiv:2402.19473 .
- [41] Zhu, X., Xie, Y., Liu, Y., Li, Y., Hu, W., 2025a. Knowledge graph-guided retrieval augmented generation. arXiv preprint arXiv:2502.06864 .
- [42] Zhu, Z., Zhao, Q., Ge, Y., Li, J., Wang, S., Yang, J.J., 2025b. Legn: A large language model-guided event graph network for intraoperative hypotension prediction. *Expert Systems with Applications* , 128677.