

Mixture of Structural-and-Textual Retrieval over Text-rich Graph Knowledge Bases

Yongjia Lei¹, Haoyu Han², Ryan A. Rossi³, Franck Dernoncourt³,
Nedim Lipka³, Mahantesh M Halappanavar⁴, Jiliang Tang², Yu Wang¹

¹University of Oregon, ²Michigan State University,

³Adobe Research, ⁴Pacific Northwest National Laboratory

{yuwang, yongjia}@uoregon.edu, {hanhaoy1, tangjili}@msu.edu

{ryrossi, dernonco, lipka}@adobe.com, hala@pnnl.gov

Abstract

Text-rich Graph Knowledge Bases (TG-KBs) have become increasingly crucial for answering queries by providing textual and structural knowledge. However, current retrieval methods often retrieve these two types of knowledge in isolation without considering their mutual reinforcement and some hybrid methods even bypass structural retrieval entirely after neighboring aggregation. To fill in this gap, we propose a Mixture of Structural-and-Textual Retrieval (MoR) to retrieve these two types of knowledge via a Planning-Reasoning-Organizing framework. In the Planning stage, MoR generates textual planning graphs delineating the logic for answering queries. Following planning graphs, in the Reasoning stage, MoR interweaves structural traversal and textual matching to obtain candidates from TG-KBs. In the Organizing stage, MoR further ranks fetched candidates by their structural trajectory. Extensive experiments demonstrate the superiority of MoR in harmonizing structural and textual retrieval with discovered insights, including uneven retrieving performance across different query logics and the benefits of integrating structural trajectories for candidate reranking. Our code is available at <https://github.com/Yoega/MoR>.

1 Introduction

Text-rich Graph Knowledge Bases (TG-KBs), due to their structured representation of textual documents, ubiquitously store textual and structural knowledge (Jin et al., 2024b). For example, scholars retrieve relevant research from paper management systems to advance scientific discoveries where nodes represent papers and edges denote references (Lu et al., 2024). With large language models (LLMs)-powered generators approaching human intelligence in language comprehension and generation, retrieving supporting knowledge from TG-KBs to contextualize and ground generation has become increasingly crucial for correctly answering queries (Gao et al., 2023b; Ni et al., 2025).

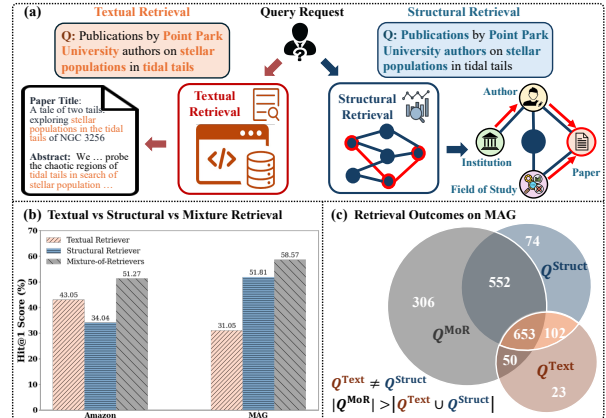


Figure 1: (a) Textual retrieval and structural retrieval. (b) The effectiveness of retrieval methods varies across different TG-KBs. (c) Within the same TG-KB, queries addressed by textual (i.e., Q^{Text}) and structural retrieval (i.e., Q^{Struct}) exhibit both overlaps and distinctiveness.

Since supporting knowledge in TG-KBs typically exhibits in both the textual and structural formats (Kolomiyets and Moens, 2011; Jin et al., 2024b), retrieval methods should be tailored to both formats effectively as shown in Figure 1(a). Textual retrieval methods retrieve textual knowledge such as indexed documents based on its similarity to the given query and can be broadly categorized into lexical methods (e.g., BM25) and semantic methods (e.g., Contriever) (Karpukhin et al., 2020; Izacard et al., 2022). Structural retrieval methods retrieve structural knowledge such as neighboring entities (Jiang et al., 2023; Edge et al., 2024; Wang et al., 2024) by applying graph traversal and graph machine learning (Yasunaga et al., 2021; Tian et al., 2024). Despite the advancements in both textual and structural retrieval, they are often applied independently and fail to mutually reinforce each other. As shown by Figure 1(b), neither structural retrieval by following the logical structure of the query nor textual retrieval by conducting Top-K BM25 matching can achieve better performance on both Amazon and MAG datasets simultaneously.

To effectively retrieve both textual and structural knowledge from TG-KBs, recent works (Xia et al., 2024; Li et al., 2024) aggregate neighboring documents to fuse structural knowledge into textual narratives, followed by textual retrieval, with Xia et al. (2024) filtering irrelevant neighbors by their relations and Li et al. (2024) weighted aggregating neighbors based on their fields. However, three challenges remain. First, rewording aggregated neighbors requires frequently invoking LLMs, resulting in prohibitive resources for long documents with exponentially growing neighbors. Second, structural signals humans use to form logical plans are completely discarded after neighbor aggregation. Third, rigid neighbor aggregation overlooks varying desires for structural and textual knowledge for different queries and TG-KBs. Even within the same TG-KB, such as MAG in Figure 1(c), queries answered by textual retrieval (i.e., Q^{Text}) are different from those by structural retrieval (i.e., Q^{Struct}).

To address the above three challenges, we infuse the mixture-of-expert philosophy into retrieval design and propose a Mixture Of Structural-and-Textual Retrieval (MoR) in Figure 2. MoR begins with a planning module that generates planning graphs to outline query logics and preserve structural signals without rewording aggregated neighbors, overcoming the first and second challenges. Next, MoR interleaves structural traversal and textual matching in the reasoning module, enabling these two retrievals to reinforce each other. Finally, MoR devises a structure-aware reranker in the organization module that adaptively adjusts the importance of retrieved textual/structural knowledge, addressing the third challenge. Via Planning–Reasoning–Organizing, MoR intelligently retrieves structural and textual knowledge based on query logical structure. Our key contributions are:

- **Planning via Textual Graph Generation:** We define retrieval planning as generating textual graphs that outline the logical structure, i.e., the plan, for identifying entities relevant to the query.
- **Reasoning via Mixture of Structural-and-Textual Traversal:** We devise a mixed traversal by interweaving textual matching and structural traversal to retrieve knowledge following query logical structure depicted by the generated plan.
- **Organizing via Structure-aware Rerank:** With candidates obtained from mixed traversal, we propose a Structure-aware Rerank to select Top-K candidates based on their traversal trajectory.

2 Preliminary

Notations: A Text-rich Graph Knowledge Base (TG-KB) \mathcal{B} generally consists of a set of connected nodes \mathcal{V} in the graph with each node $v \in \mathcal{V}$ associated with its corresponding document $\mathcal{D}_v \in \mathcal{D}$ and category $\mathcal{E}_v \in \mathcal{E}$. When retrieving nodes with supporting documents from \mathcal{B} for answering a given query $Q \in \mathcal{Q}$, we typically follow certain rationale encapsulating the underlying logic of that query (Xu et al., 2024; Xue et al., 2024), which can be characterized by a text-attributed planning graph G . In many existing works (Jin et al., 2024a; Wu et al.), this planning graph can be usually decomposed into multiple reasoning paths $G = \{\mathcal{P}_i\}_{i=1}^{|\mathcal{G}|}$ where the i^{th} reasoning path $\mathcal{P}_i = (p_{i1} \rightarrow p_{i2} \rightarrow \dots \rightarrow p_{iL_i})$ is a distinctive reasoning chain of length L_i encoding a unique logic and the j^{th} node p_{ij} corresponds to an entity in \mathcal{B} with its own category $\mathcal{E}_{p_{ij}}$ and textual restriction $\mathcal{T}_{p_{ij}}$ extracted from the query. For example, in Figure 1(a), the query *Publications by Point...* has a planning graph with two paths, i.e., $\mathcal{P}_1 = (\text{Institution} \rightarrow \text{Author} \rightarrow \text{Paper})$ and $\mathcal{P}_2 = (\text{Field-of-Study} \rightarrow \text{Paper})$, where the category and textual restriction of the first node on \mathcal{P}_1 are $\mathcal{E}_{p_{11}} = \text{Institution}$ and $\mathcal{T}_{p_{11}} = \langle \text{Point Park Univerisity} \rangle$, respectively. Comprehensive notations are summarized in Table 6 in Appendix A.

Problem Setup: With the above notations, the investigated problem is to retrieve entities $\mathcal{C} \subseteq \mathcal{V}$ answering a given query Q .

Textual Retrieval retrieves candidates based on the textual signals of both the query and documents. One common strategy is to retrieve candidates $\tilde{\mathcal{C}}$ from the whole documents \mathcal{D} that have Top-K textual similarity to query Q measured by lexical or semantic similarity (Vijaymeena and Kavitha, 2016). The textual retrieval used in MoR retrieves documents for a given query by matching them with textual descriptions in the query, e.g., matching *stellar populations in tidal tails* shown in Figure 1.

Structural Retrieval retrieves candidates by applying prescribed rules to structured databases such as knowledge graphs and SQL (Guo et al., 2023). Common strategies include graph-based traversal (e.g., BFS, DFS) and rule fetching (Jiang et al., 2023). Specifically, MoR conducts structural retrieval by traversing neighbors of certain categories from the generated planning graph. For example, in Figure 1(a), only "Paper" typed neighbors of the Author can be traversed by our structural retrieval.

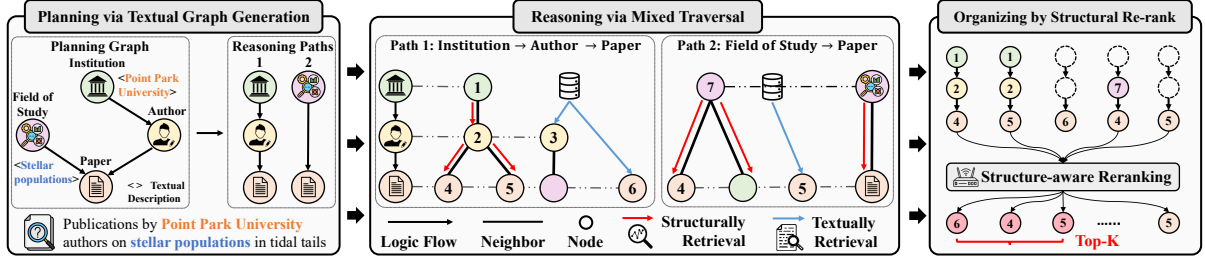


Figure 2: Our MoR framework consists of a planning module to generate a planning graph, a reasoning module to conduct mixed traversal, and an organizing module to rerank the retrieved candidates.

3 Framework

In a nutshell, we formulate our MoR as the conditional distribution $P_{\Theta}(\mathcal{C}|Q, \mathcal{B})$ of retrieved candidates \mathcal{C} given the user input query Q over TG-KB \mathcal{B} , which is further factorized into three distributions corresponding to our proposed three modules: planning via generating the text-attributed planning graph G , reasoning via conducting mixture of structural-and-textual traversal to obtain intermediate candidates $\tilde{\mathcal{C}}$ following the generated planning graph G , and organizing via applying structure-aware reranking to the obtained candidates $\tilde{\mathcal{C}}$, obtaining final candidates \mathcal{C} :

$$P_{\Theta}(\mathcal{C}|Q, \mathcal{B}) = \sum_{G \in \mathbb{G}} \left[\sum_{\tilde{\mathcal{C}} \in \mathbb{C}} \underbrace{P_{\Theta_3}(\tilde{\mathcal{C}}|G, Q, \mathcal{B})}_{\text{Organizing}} \right] \times \underbrace{P_{\Theta_2}(\tilde{\mathcal{C}}|G, Q, \mathcal{B})}_{\text{Reasoning}} \times \underbrace{P_{\Theta_1}(G|Q, \mathcal{B})}_{\text{Planning}}$$

where $P_{\Theta_1}(G|Q, \mathcal{B})$ is the probability distribution of generating the text-attributed planning graph G given the input query Q and TG-KB \mathcal{B} ; $P_{\Theta_2}(\tilde{\mathcal{C}}|G, Q, \mathcal{B})$ is the probability distribution of retrieving intermediate candidates $\tilde{\mathcal{C}}$ given the planning graph G and the query Q via our mixed traversal; $P_{\Theta_3}(\mathcal{C}|\tilde{\mathcal{C}}, G, Q, \mathcal{B})$ is the probability distribution of reranking the intermediate candidates so that Top-K positions form the ground-truth entities \mathcal{C} . \mathbb{G}/\mathbb{C} denotes the space of all possible planning graphs and all possible configurations of size-K candidate node sets from all nodes \mathcal{V} in TG-KB \mathcal{B} . The overall objective is to maximize the likelihood of retrieving ground-truth candidates \mathcal{C} for each input query $Q \in \mathcal{Q}$:

$$\Theta^* = \arg \max_{\Theta} \prod_{Q \in \mathcal{Q}} P_{\Theta}(\mathcal{C}|Q, \mathcal{B}) \quad (1)$$

Following the above paradigm, we next introduce the three components: Planning via textual graph generation in Section 3.1, Reasoning via mixed traversal in Section 3.2, and Organizing via structure-aware reranking in Section 3.3.

3.1 Planning via Textual Graph Generation

To effectively reason over the underlying logic of queries and answer them, we propose a planning module that constructs a planning graph to capture their underlying logical structures. Unlike conventional approaches relying on rigid heuristics, e.g., shortest-path retrieval in knowledge graphs (Luo et al.; Delile et al., 2024), or step-by-step prompting LLMs, which incur high computational costs (Sun et al., 2023; Wang et al., 2024), our method generates the entire planning graph in one shot, eliminating repeated LLM calls. More importantly, as planning graphs integrate entity restrictions encoding query-specific constraints and entity categories capturing broader logical structure, our MoR can generalize learned patterns and efficiently adapt to new queries with the same underlying logic. For example, any query with the form *Papers associated with <institution> and are in the field of <field>* shares the same patterns with the query in Figure 2. Below, we first formalize the planning graph and then optimize its generation.

3.1.1 Planning Graph Formulation

A planning graph G is a structured representation where nodes represent entities and edges denote their logical relations. Each entity is associated with both a category and query-specific restriction. For example, given the query *Can you give me publications by Point Park University authors on stellar populations in tidal tails*, the generated planning graph is: $G = (\text{Institution} \langle \text{Point Park University} \rangle \rightarrow \text{Author} \rightarrow \text{Paper} \leftarrow \text{Field-of-Study} \langle \text{Stellar Population} \rangle)$ with Institution, Author, Paper, Field-of-Study as categories and $\langle \text{Point Park University} \rangle$, $\langle \text{Stellar Populations} \rangle$ as restrictions. Note that edges in our planning graph can also possess different categories. For example, in the biomedical TG-KBs, the relation between Disease and Drug entities could be Indication or Contra-indication (Wu et al.), adding a finer level of semantic distinction to the relation.

3.1.2 Planning Graph Optimization

To ensure that our generated planning graph captures the query logic, we train a textual graph generator to maximize the likelihood of generating ground-truth planning graphs given their queries. Formally, given the joint distribution of the training pairs between queries and planning graphs $P_{\mathcal{Q} \times \mathcal{G}}^{\text{Train}}$, we optimize the planning module P_{Θ_1} by solving:

$$\arg \max_{\Theta_1} \mathbb{E}_{(Q,G) \sim P_{\mathcal{Q} \times \mathcal{G}}^{\text{Train}}} \log P_{\Theta_1}(G|Q, \mathcal{B}) \quad (2)$$

To avoid the combinatorial explosion of exponentially growing planning graph candidates (You et al.), we decompose each planning graph into multiple reasoning paths $G = \{\mathcal{P}_i\}_{i=1}^{|G|}$. Each path $\mathcal{P}_i = (p_{i1} \rightarrow, \dots, \rightarrow p_{iL_i})$ represents a distinct reasoning chain, where node p_{ij} denotes an entity in TG-KB sharing the same textual category $\mathcal{E}_{p_{ij}}$ and satisfying the restriction $\mathcal{T}_{p_{ij}}$ from the query. Given the sequential nature and textual formats of these decomposed reasoning paths, LLMs can be naturally employed here as the planning graph generator, which conducts next-token prediction by predicting j^{th} token t_j conditioned on preceding tokens $t_{<j}$, the query Q and the TG-KB \mathcal{B} :

$$P_{\Theta_1}(G|Q) = \prod_{j=1}^n P_{\Theta_1}(t_j | t_{<j}, Q, \mathcal{B}). \quad (3)$$

Note that our proposed planning graph generator is not limited to LLMs. Any graph generative model preserving both structural dependencies and textual associations can be employed (Zhu et al.).

3.2 Reasoning via Mixed Traversal

Following the reasoning paths of the above planning graph $G = \{\mathcal{P}_i\}_{i=1}^{|G|}$, the reasoning module conducts a mixed traversal by interweaving neighbor fetching and textual matching to form intermediate candidates $\tilde{\mathcal{C}}$, which are introduced next.

3.2.1 Structural Traversal

Following the definition in Section 2 that structural retrieval follows prescribed rules for knowledge retrieval, here we set these prescribed rules to be iteratively performing layer-wise breadth-first-search that traverses neighboring entities with categories aligning with those in the reasoning paths. Concretely, reasoning at the l^{th} -step of the planning path \mathcal{P}_i , we check for each node v in candidates set of last layer $\forall v \in \tilde{\mathcal{C}}_{i-1}^l$ and fetch its neighbors

$\forall u \in \mathcal{N}_v$ with the same category as the corresponding node p_{il} (i.e., $\mathcal{E}_u = \mathcal{E}_{p_{il}}$) in the reasoning path, which can be mathematically formulated as:

$$\tilde{\mathcal{C}}_i^{l, \text{Struct}} = \cup_{v \in \tilde{\mathcal{C}}_{i-1}^l} \{u | u \in \mathcal{N}_v, \mathcal{E}_u = \mathcal{E}_{p_{il}}\} \quad (4)$$

where $\tilde{\mathcal{C}}_i^{l, \text{Struct}}$ denotes the set of structurally retrieved entities at the l^{th} reasoning step according to the path \mathcal{P}_i and $\mathcal{E}_u = \mathcal{E}_{p_{il}}$ ensures that the category of the traversed neighbor u matches the corresponding entity category routine by the planning graph, resonating the nature of rule-based structural retrieval. Note that the seeding candidates $\tilde{\mathcal{C}}_i^{1, \text{Struct}}$ at the very first layer are initialized by retrieving Top-K entities through textual matching, i.e., $\tilde{\mathcal{C}}_i^{1, \text{Struct}} = \tilde{\mathcal{C}}_i^{1, \text{Text}}$, which is introduced next.

3.2.2 Textual Matching

In addition to retrieving structural knowledge, our MoR also retrieves textual knowledge via Textual Matching, which retrieves candidates based on their textual similarity to queries. For each reasoning node p_{il} at l^{th} reasoning step along the reasoning path \mathcal{P}_i , we concatenate the query and the textual restriction of p_{il} , i.e., $Q' = [Q : \mathcal{T}_{p_{il}}]$, then compute its textual similarity to documents of nodes in TG-KB, i.e., $\phi(Q', \mathcal{D}_v), \forall v \in \mathcal{V}$, and finally retrieve the Top-K scored candidates:

$$\tilde{\mathcal{C}}_i^{l, \text{Text}} = \text{TopK}(\{v | v \in \mathcal{V}, \mathcal{E}_v = \mathcal{E}_{p_{il}}, \phi(Q', \mathcal{D}_v)\}) \quad (5)$$

Integrating candidates from structural traversal and textual matching together, the final candidates at l^{th} -step of \mathcal{P}_i are formed as:

$$\tilde{\mathcal{C}}_i^l = \tilde{\mathcal{C}}_i^{l, \text{Struct}} \cup \tilde{\mathcal{C}}_i^{l, \text{Text}}, \forall l \in \{1, 2, \dots, L_i\} \quad (6)$$

The integrated candidates $\tilde{\mathcal{C}}_i^l$ serve as seeding nodes initializing the next round of planning graph-guided structural traversal and textual matching, which creates a mutual reinforcement between structural and textual knowledge since previously retrieved two knowledge can both inform next round of structural/textual knowledge retrieval.

We iteratively conduct mixed traversal for every reasoning path $\mathcal{P}_i \in G$ and integrate retrieved entities together by taking their intersection, i.e., $\tilde{\mathcal{C}} = \cap_{\mathcal{P}_i \in G} \tilde{\mathcal{C}}_i^{L_i}$, adhering to the fact that candidates should simultaneously satisfy the logic routine by all reasoning paths. Note that no training is involved in the mixed graph traversal, i.e., $P_{\Theta_2}(\tilde{\mathcal{C}}|G, Q, \mathcal{B}) = P(\tilde{\mathcal{C}}|G, Q, \mathcal{B})$. Future works could explore optimizing graph traversal by rewards from agent-environment interactions (Nguyen et al., 2024).

3.3 Organizing via Structure-aware Rerank

Although the retrieved candidates from Section 3.2 strictly adhere to the prescribed rule given by the planning graph, the sheer volume of candidates misaligns with realistic constraints (e.g., Top-20 retrieval budget (Zeng et al., 2024)) and may even cause difficulty to downstream executors such as long-context challenges for LLMs. To better emulate human reasoning, where multiple clues are gathered, analyzed in relation to the query, and synthesized into a coherent answer, we propose a structure-aware reranker to organize and rerank the candidates $\tilde{\mathcal{C}}$, and select Top-K ones as the final retrieved answers \mathcal{C} . Instead of relying only on textual features (Hu et al., 2019), our reranker assigns a ranking score based on features of structural trajectories obtained from the mixed traversal in Section 3.2, innovatively leveraging both structural and textual knowledge in reranking.

Previously, $\tilde{\mathcal{C}}$ is defined as intermediate retrieved entities. To consider structural features in reranking, we pair each retrieved candidate in $\tilde{\mathcal{C}}$ with its corresponding traversal trajectory obtained from the reasoning module. Specifically, each trajectory \mathcal{P}_i of length L_i is featuring three types of attributes:

- **Textual Fingerprint (TF)**: Concatenation of similarity scores between the expanded query and each node on the path: $\|_{l=1}^{L_i} \phi(Q', \mathcal{D}_{p_{il}})$.
- **Structural Fingerprint (SF)**: Concatenation of node categories at each step on the path: $\|_{l=1}^{L_i} \mathcal{E}_{p_{il}}$
- **Traversal Identifier (TI)**: Concatenation of the indicator specifying whether each step uses a structural or textual retrieval: $\|_{l=1}^{L_i} \mathcal{T}_{p_{il}}$.

We then train a reranker on these trajectories using the cross-entropy loss. For a training query Q and its associated candidate trajectory \mathcal{P}_i , the loss is computed as follows:

$$\mathcal{L}_{\Theta_3} = - \sum_{(\mathcal{P}_i, Q) \in \tilde{\mathcal{C}}} \sum_{j=1}^2 y_j^i \log(\sigma(f(\underbrace{\|_{l=1}^{L_i} \mathcal{E}_{p_{il}}}_{\text{Structural Fingerprint}} : \underbrace{\|_{l=1}^{L_i} \phi(Q', \mathcal{D}_{p_{il}})}_{\text{Textual Fingerprint}} : \underbrace{\|_{l=1}^{L_i} \mathcal{T}_{p_{il}}}_{\text{Traversal Identifier}}))_j). \quad (7)$$

where $f(\cdot)$ is the reranker producing a score for each (Q, \mathcal{P}_i) pair, $\sigma(\cdot)$ denotes the softmax function, and $y_j^i \in \{0, 1\}$ indicates whether the i -th candidate is a correct (positive) or incorrect (negative) match for Q . This formulation encourages the reranker to assign higher scores to positive trajectories, thereby improving ranking performance.

4 Experiment

4.1 Experimental Setup

We briefly introduce experimental settings to verify our proposed MoR, including Datasets & Baselines, Implementation Details, and Evaluation Metrics. More details are in Appendix B.

Datasets & Baselines: We use three TG-KBs from STaRK (Wu et al.) covering three knowledge domains, including E-commerce Products (Amazon), Academic Papers (MAG), and Biomedicine (Prime). We compare our MoR with baselines established by Wu et al. and categorize them into textual/structural/hybrid-based ones. More recent state-of-the-art hybrid retrieval approaches from TG-KBs such as KAR (Xia et al., 2024) and MFAR* (Li et al., 2024) are also compared.

Implementation Details: To enhance the planning capability of our planning module, we fine-tune the Llama 3.2 (3B) on 1000 sampled queries with their corresponding ground-truth planning graphs, serving as the textual graph generator. In the absence of ground-truths, we synthesize them using LLMs. For the Prime dataset, we empirically find that directly prompting LLMs can hardly generate accurate planning graphs due to the lack of biomedical domain knowledge (Shen et al.). Therefore, we adopt an alternative approach. First, we instruct LLMs to extract triplets from each query and then construct the planning graphs by merging triplets with shared entities. During mixed traversal, textual matching can be implemented using any lexical or semantic methods. For this study, we employ BM25 for Amazon and MAG and fine-tune a retriever to complement the biomedical knowledge for Prime. To initialize the structural traversal, we employ textual matching to locate the top 5 nodes that are most relevant to the query as seeds. Additionally, at each layer, we incorporate the top 10 nodes retrieved via textual matching and append them to the current candidate set for the next round of traversal. Notably, due to the uncertainty of LLMs, the generated planning graphs can be invalid. In this case, we will directly conduct textual matching to retrieve candidates. For our ablations without reranker, we employ Ada-002 (Wu et al.) with cosine similarity as the scorer to rank candidates for evaluating performance.

Evaluation Metrics: We follow Wu et al. for evaluation by reporting Hit@1 (H@1), Hit@5 (H@5), Recall@20 (R@20), and mean reciprocal rank MRR to evaluate in the full spectrum.

Category	Retrieval Baseline	AMAZON				MAG				PRIME				AVERAGE			
		H@1	H@5	R@20	MRR	H@1	H@5	R@20	MRR	H@1	H@5	R@20	MRR	H@1	H@5	R@20	MRR
Textual	BM25 (Wu et al.)	44.94	67.42	53.77	55.30	25.85	45.25	45.69	34.91	12.75	27.92	31.25	19.84	27.85	46.86	43.57	36.68
	Ada-002 (Wu et al.)	39.16	62.73	53.29	50.35	29.08	49.61	48.36	38.62	12.63	31.49	36.00	21.41	26.96	47.94	45.88	36.79
	Multi-ada-002 (Wu et al.)	40.07	64.98	55.12	51.55	25.92	50.43	50.80	36.94	15.10	33.56	38.05	23.49	27.03	49.66	47.99	37.33
	DPR (Karpukhin et al., 2020)	15.29	47.93	44.49	30.20	10.51	35.23	42.11	21.34	4.46	21.85	30.13	12.38	10.09	35.00	38.91	21.31
Structural (KG)	QAGNN (Yasunaga et al., 2021)	26.56	50.01	52.05	37.75	12.88	39.01	46.97	29.12	8.85	21.35	29.63	14.73	16.10	36.79	42.88	27.20
	ToG (Sun et al., 2023)	-	-	-	-	13.16	16.17	11.30	14.18	6.07	15.71	13.07	10.17	9.62	15.94	12.18	12.18
Hybrid	AvaTaR (Wu et al., 2025)	49.87	69.16	60.57	58.70	44.36	59.66	50.63	51.15	18.44	36.73	39.31	26.73	37.56	55.18	50.17	45.53
	KAR (Xia et al., 2024)	54.20	68.70	57.24	<u>61.29</u>	50.47	69.57	60.28	58.65	30.35	49.30	50.81	39.22	45.01	62.52	56.11	53.05
	MFAR* (Li et al., 2024)	41.20	<u>70.00</u>	58.50	54.20	49.00	<u>69.60</u>	<u>71.70</u>	58.20	40.90	62.80	68.30	51.20	43.70	<u>67.47</u>	66.17	<u>54.53</u>
	HYBGRAG (Lee et al., 2024)	-	-	-	-	65.40	75.31	65.70	69.80	28.56	41.38	43.58	34.49	50.91	58.35	54.64	52.15
	MoR	<u>52.19</u>	74.65	<u>59.92</u>	62.24	<u>58.19</u>	78.34	75.01	<u>67.14</u>	<u>36.41</u>	<u>60.01</u>	<u>63.48</u>	<u>46.92</u>	<u>48.93</u>	71.00	<u>66.14</u>	58.77
Ablation	MoR _{w/o R}	44.21	68.87	56.50	55.28	34.33	62.55	67.55	47.40	31.59	53.48	60.74	41.81	31.07	57.04	57.73	43.03
	MoR _{w/o RT}	34.04	53.41	45.16	42.85	51.81	73.54	74.17	61.68	28.95	46.12	49.54	36.56	36.39	56.73	55.73	45.53
	MoR _{w/o RS}	43.05	69.36	57.38	54.69	31.05	51.84	50.56	40.64	22.27	38.45	39.21	29.41	28.95	51.28	48.02	38.98

Table 1: Comparing different retrieval methods with our proposed MoR and its ablations on Amazon, MAG, and Prime datasets. The best and runner-up results are in **bold** and underlined. Overall, MoR achieves the best performance. Note that MFAR* denotes the best model variant proposed in (Li et al., 2024)

4.2 Overall Retrieval Performance

We compare MoR with other baselines on three TG-KBs in Table 1. Generally, hybrid methods, AvaTaR, KAR, MFAR*, and our MoR, achieve better performance than purely textual or structural methods owing to their ability to integrate both structural and textual knowledge. Among all baselines, our proposed MoR achieves the overall best performance with a substantial margin on average, with the first ranking on MAG and the second ranking on Amazon/Prime datasets. This demonstrates the effectiveness of our proposed mixture of structural and textual knowledge retrieval. Textual retrieval performs better on Amazon than on MAG, suggesting that Amazon queries rely more on textual knowledge. In contrast, its weaker performance on MAG is due to MAG’s lower textual richness and stronger structural signals. This disparity aligns with the distribution analysis presented by Wu et al. and supports our hypothesis that queries in different TG-KB datasets require varying desires for textual and structural knowledge. Meanwhile, structural retrieval methods such as conventional knowledge graph-based ones perform poorly because they are designed for graphs with minimal textual information compared to TG-KBs. Different from Amazon and MAG, all existing methods without supervised tuning (e.g., Ada-002) exhibit significantly lower performance on Prime. This is due to the extreme domain expertise required in biology, where word-count-based, pre-trained textual similarity-based, and even more powerful LLMs are all poorly applicable here. Through fine-tuning, MFAR* and our proposed MoR generally achieve better performance, demonstrating the necessity of domain-specific knowledge for answering queries in knowledge-intensive domains.

4.3 Ablation Study

After verifying the superiority of MoR, we conduct ablation studies to assess its different components, including module and feature ablation.

4.3.1 Module Ablation

To assess the contribution of each module in MoR, namely, Text Matching-based Retrieval, Neighborhood-Fetching-based Structural Retrieval, and Reranker, we conduct a series of ablation experiments. First, we remove the Reranker, resulting in the variant MoR_{w/o R}. On top of that, we further separately eliminate Text Retrieval and Structural Retrieval, yielding MoR_{w/o RT} and MoR_{w/o RS}, respectively. As shown in Table 1, the complete MoR framework consistently achieves the highest performance across all datasets, demonstrating the synergistic effect of the Textual Retriever, Structural Retriever, and Reranker. After removing Reranker, MoR_{w/o R} exhibits a consistent performance drop across all datasets and evaluation metrics. This underscores the importance of the Reranker in refining retrieval by filtering noisy candidates from the intermediate reasoning stage. Eliminating Text Retrieval, i.e., MoR_{w/o RT}, leads to a notable performance drop on Amazon but an unexpected improvement on MAG. This suggests that while textual knowledge benefits Amazon, it introduces misleading hard negatives that compromise the ranking method (e.g., Ada-002) for MAG. Conversely, removing Structural Retrieval, MoR_{w/o RS}, results in a slight performance decrease further on MAG, reinforcing the importance of structural knowledge in MAG-related queries. These results underscore the Reranker’s crucial role in adaptively harmonizing, balancing, and selecting knowledge from both structural and textual retrieval experts.

Dataset	TF	SF	TI	H@1	H@5	R@20	MRR
MAG	✓	✗	✗	48.96	73.02	72.44	59.79
	✗	✓	✗	18.79	41.91	52.85	29.84
	✗	✗	✓	18.16	41.53	52.78	29.31
	✓	✓	✗	58.04	77.14	74.42	66.75
	✓	✗	✓	58.16	77.59	74.96	66.85
	✗	✓	✓	17.93	38.01	46.79	27.48
	✓	✓	✓	58.19	78.34	75.01	67.14
Amazon	✓	✗	✗	51.21	74.05	59.79	61.27
	✗	✓	✗	8.09	24.48	25.62	16.94
	✗	✗	✓	5.84	16.62	12.94	11.57
	✓	✓	✗	50.91	73.38	59.58	61.15
	✓	✗	✓	51.09	73.56	59.61	61.14
	✗	✓	✓	8.09	24.48	25.62	16.94
	✓	✓	✓	52.19	74.65	59.92	62.24
Prime	✓	✗	✗	35.23	59.44	63.15	46.02
	✗	✓	✗	12.95	30.48	43.33	21.44
	✗	✗	✓	11.81	27.81	40.36	19.73
	✓	✗	✓	35.80	60.12	63.40	46.50
	✓	✓	✗	35.70	60.01	63.21	46.19
	✗	✓	✓	13.20	32.27	48.01	22.95
	✓	✓	✓	36.41	60.01	63.48	46.92

Table 2: Ablation study investigating the importance of three features, Textual Fingerprint (TF), Structural Fingerprint (SF), and Traversal Identifier (TI), of the traversal trajectories used in our Structure-aware Reranker.

4.3.2 Feature Ablation

The above ablation study highlights the crucial role of Structure-aware Reranker in adaptively integrating structural and textual knowledge. To further analyze the contributions of its three key features, **Textual Fingerprint (TF)**, **Structural Fingerprint (SF)**, and **Traversal Identifier (TI)** defined in Section 3.3, we conduct a feature ablation analysis and report retrieval performance across different feature configurations in Table 2. Overall, using three features together yields the best performance on both MAG and Amazon, highlighting their synergistic effect. Individually, TF contributes the most and outperforms SF and TI on both datasets. The reason is that based on the definition in Section 3.3, TF directly captures the relevance between the query and the retrieved nodes along the trajectory, whereas SF and TI primarily characterize the structural patterns and retrieval types, serving more as complementary factors. Therefore, equipping TF with these complementary factors (i.e., SF or TI) yields around 10% additional gains on MAG. This is because SF and TI help the reranker selectively emphasize the relevance scores given by TF for certain nodes along the path. However, this boost is not observed on Amazon. We hypothesize that the textual knowledge needed there is predominantly derived from the final node on each path, making the structural cues provided by SF and TI less beneficial and even prone to overfitting. A deeper analysis to further justify this hypothesis is in Section 4.4. Overall, these findings underscore the varying importance of structural features in ranking across datasets.

Feature	MAG			Amazon			Prime		
	H@1	R@20	MRR	H@1	R@20	MRR	H@1	R@20	MRR
Last Node	49.91	73.49	59.92	50.36	59.62	61.05	33.52	61.95	44.15
Full Path	58.19	75.01	67.14	52.19	59.92	62.24	36.41	63.48	46.92

Table 3: Comparing reranking performance using last node in the retrieved trajectory and the whole trajectory.

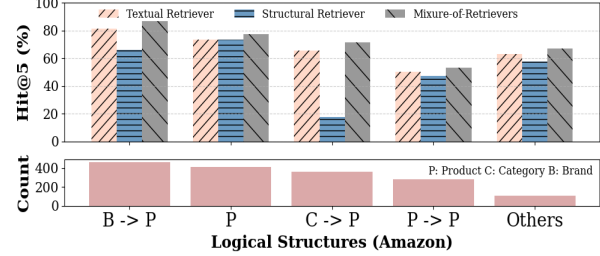


Figure 3: Imbalance number of queries and performance of different retrievers across different logical structures.

4.4 Further Analysis

This section understands MoR’s behavior by examining three questions, each of which enriches our insight into MoR’s functionality and offers novel perspectives inspiring future query retrieval research. More analysis can be found in Appendix C.

Do structure signals affect reranking? To assess the impact of trajectory information on the Reranker’s decision-making, we introduce a node-based Reranker that constructs trajectory features using only TF/SF/TI of the last node. In Table 3, the path-based Reranker outperforms the node-based variant, especially on MAG. This highlights the critical role of trajectory features/structural knowledge in reranking. The minor performance boost on Amazon after switching to the full path trajectory indicates its textual knowledge preference over the last node rather than the whole trajectory.

How does MoR perform on different logical structures? Figure 3 shows the average performance of MoR on each query group categorized by their logical structures, where "Others" refer to queries with undefined logical structures in Wu et al. MoR consistently outperforms structural and textual retrievers across different logical structures. Among all queries, MoR performs the worst on "P → P" queries due to the ambiguity uniquely caused by repeated "product" entities from multi-step traversal. The average-performing "Others" group underscores the utility of diverse planning strategies for the same query. Lastly, the skewed query distribution and retrieval performance across planning patterns reflect the varying nature of real-world planning needs. We hope these insights inspire data-centric reasoning and error control of planning for heterogeneous query structures.

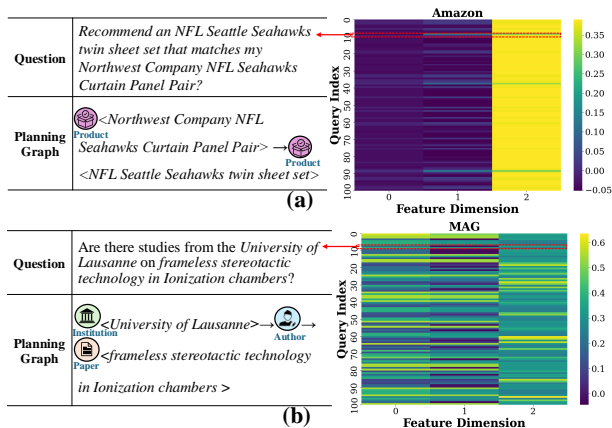


Figure 4: Saliency map visualization of query attention over three entities along the retrieved paths.

Does MoR indeed adaptively leverage the trajectory knowledge? To understand how our proposed reranker prioritizes candidates in the Top-K results, we visualize the saliency map by computing the gradient of ranking scores with respect to the textual fingerprint (TF) of three nodes along the traversed path, which quantifies their importance for answering a given query. Figure 4 illustrates this by analyzing trajectories for 100 ground-truth candidates across 100 queries on the Amazon and MAG datasets. Each dimension corresponds to a traversed node, with the final one representing the candidate itself. While the saliency score is concentrated in the last dimension for Amazon, MAG exhibits a more evenly distributed saliency pattern, where multiple nodes along the path contribute significantly to ranking score computation. This suggests that structural knowledge is more critical for answering queries in MAG, aligning with the previously observed lower performance of purely textual retrieval on MAG in Table 1. Further case studies explain why the reranker attends different nodes for different queries. In Figure 4(a), the reranker favors the last two dimensions as the rich textual restriction (i.e., "Northwest Company..." and "NFL Seattle...") aids in identifying the correct node at the corresponding reasoning step, as discussed in Section 3.2. These correct nodes with higher similarity scores with the query help guide the retrieval process toward the ground truth. Conversely, in Figure 4(b), since the first node ("University of Lausanne") helps narrow the search space and the last node ("frameless...") further filter candidates, both nodes have high saliency scores. Overall, our findings demonstrate that the reranker dynamically adapts its reliance on structural and textual knowledge depending on the dataset and query.

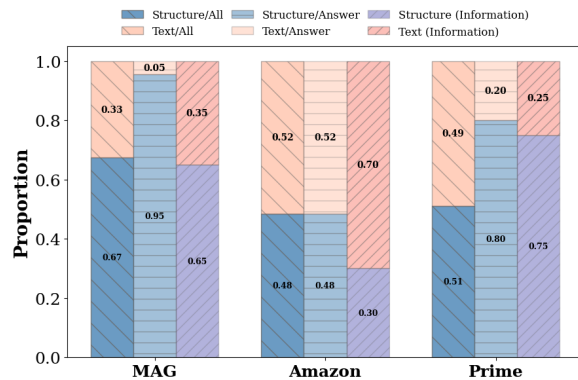


Figure 5: Traversal Identifier analysis in Top-20 retrieved candidates and information ratio in dataset.

Does MoR adaptively prioritize retrieval methods for different datasets? To assess the degree to which MoR adaptively leverages structural and textual retrieval across datasets, we calculate the following three ratios: **Structure/All**, **Text/All**: Fraction of structurally or textually retrieved candidates within the Top-20 retrieved candidates; **Structure/Answer**, **Text/Answer**: Fraction of ground-truth candidates that are structurally or textually retrieved among the Top-20; **Structure (Information)**, **Text (Information)**: Ratio of word counts of sampled relation/property requirements used in query construction.

The results are shown in Figure 5. On the MAG dataset, a majority of the Top-20 retrieved candidates (67%) and over 95% of the correct answers originate from structural retrieval. This underscores the critical role of structural signals in this domain. The finding is consistent with the inherently rich structural nature of MAG, evidenced by the higher proportion of **Structure (Information)** compared to **Text (Information)**, and is further supported by the substantial performance gains observed when structural features are incorporated, as demonstrated in Table 2 and Table 3. For the Amazon dataset, textual matching accounts for 52% of the Top-20 candidates and contributes more answers than structural traversal. This aligns with the rich textual content observed in Amazon, as shown in Figure 5. Overall, these results demonstrate that MoR exhibits adaptive retrieval behavior across datasets, effectively prioritizing the most informative retrieval strategy based on the underlying data characteristics. This provides strong evidence that MoR can dynamically coordinate between structural and textual retrieval sources.

Method	Planning	Reasoning	Organizing
Theoretical	$\mathcal{O}(K \cdot D)$	$\mathcal{O}(Ad^{L-1})$	$\mathcal{O}(Bd^{L-1})$
Empirical (s)	0.971	0.463	0.0340
Empirical+Batch (s)	-	-	0.0134
Empirical+Parallel (s)	-	0.398	-

Table 4: Theoretical and empirical time complexity for each component during single-query retrieval. Batch and parallel optimization are applied to the Organizing and Reasoning stages, respectively.

Dataset	1-hop	2-hop	3-hop
Prime	0.059 s	0.068 s	0.072 s
Amazon	0.381 s	0.403 s	0.669 s
MAG	0.289 s	0.297 s	1.233 s

Table 5: Empirical traversal time across different hop depths for each dataset.

4.5 Efficiency and Scalability of MoR

Since MoR consists of three stages, we conduct the complexity analysis by analyzing each of these stages theoretically and empirically. The results are shown in Table 4, where K means the number of tokens in the generated planning graph, D is the model dimension following the GPT-style decoding and key-value caching, d denotes the average degree of node at each step, L means the number of steps/layers in one path, and both A and B are constant. Due to the limitations of LLMs, it is hard to improve the efficiency of **Planning**. As the traversal is independent, we have implemented the parallel version to speed up **Reasoning**. For the **Organizing** stage, we can easily use batch to compute the ranking scores for multiple candidates and queries simultaneously, improving the efficiency.

Based on the above analysis, the time complexity for the whole framework is $\mathcal{O}(K \cdot D + (A + B)d^{L-1}) \approx \mathcal{O}(d^{L-1})$, indicating that the complexity exponentially grows as the length of the reasoning path increases. To examine how time complexity varies with path length, we empirically analyze the scalability of MoR with respect to the depth of reasoning paths. We group test queries by the number of reasoning hops (e.g., 1-hop, 2-hop, 3-hop) and compare efficiency across these three groups on three datasets. As shown in Table 5, most queries can be processed within 1s. While deeper queries naturally require more processing time, MoR maintains significantly high efficiency due to its step-wise traversal and adaptive fallback to textual retrieval when necessary.

More detailed analysis of computational complexity and scalability to large datasets can be found in Appendix D.1 and Appendix D.2.

5 Related Work

Retrieval-augmented Generation (RAG): RAG enhances generative tasks by retrieving relevant information from external knowledge sources (He et al., 2025; Gao et al., 2023b) and has been widely used to improve question-answering (Liu et al., 2023). With LLMs, RAG has been used for mitigating hallucinations (Yao et al., 2023), enhancing interpretability (Gao et al., 2023a), and enabling dynamic knowledge updates (Wang et al., 2024). This work essentially leverages the idea of RAG to retrieve supporting entities from TG-KBs to contextualize answer generation. Depending on concrete types of knowledge being retrieved, existing retrievers can be categorized into structural and textual retrieval, which are reviewed next.

Textual and Structural Retrieval: Early textual retrieval models, such as TF-IDF and BM25 (Robertson et al., 2009), rely on lexical similarity and keyword matching (Chen et al., 2017; Yang et al., 2019; Mao et al., 2021). Modern approaches address this limitation by learning dense representations (Karpukhin et al., 2020). Beyond textual retrieval, structural retrieval leverages graph-based techniques to extract structured knowledge. Methods such as graph traversal (Wang et al., 2024; Jiang et al., 2023), community detection (Edge et al., 2024), and graph machine learning models, including graph neural networks (Yasunaga et al., 2021; Mavromatis and Karypis, 2024), play a crucial role in structural retrieval. Our approach integrates the strengths of both textual and structural retrieval by infusing the mixture-of-expert philosophy into retrieval design.

Due to page limitation, a comprehensive version of the related work is attached in Appendix F.

6 Conclusion

In this work, we propose a mixture of structural and textual retrieval (MoR) to adaptively retrieve structural and textual knowledge based on query desire, which first utilizes a textual graph generator to generate the planning graph, then performs a mixed traversal and conducts organizing via a structure-aware reranker to obtain final candidates. Experiments demonstrate the advantages of our MoR in harmonizing the retrieval of both textual and structural knowledge with insightful discoveries, including balancing retrieval performance across queries with different patterns and query-adaptive knowledge desire for structural/textual knowledge.

7 Limitations

In this paper, we integrate a mixture of expert philosophy into retrieval design and propose a Mixture of structural-and-textual Retrieval (MoR) to adaptively retrieve textual and structural knowledge. The limitations of MoR can be categorized into two main aspects in the following:

Lack of Domain-Specific Knowledge: Our proposed MoR, similar to other baselines, does not exhibit significantly higher performance on PRIME than AMAZON and MAG. The reason is the lack of biomedical knowledge required to comprehend biomedical questions, extract key information, navigate relevant entities and relations, and rerank retrieved candidates. This suggests that current state-of-the-art retrieval models, even paired with LLMs’ intelligence, still struggle to handle domain-specific knowledge effectively. Such limitations may extend to other specialized domains, such as finance and law. Future research could integrate domain-specific knowledge into retrieval.

Reranking at Every Traversal Layer: Our current MoR adaptively routes retrieved candidates into the Top-K positions at the final layer via reranking, effectively implementing a conventional Mixture of Experts (MoE) routing mechanism. Despite the state-of-the-art performance we have achieved in Table 1, this routing mechanism could also be applied to intermediate layers, where after each retrieval step, candidates are reranked, and only Top-K proceeds to the next round of traversal and retrieval. This enables every layer of mixed traversal to emulate the router design of the MoE.

Multi-Trajectory Reranking: While our current Structural Reranker is designed to compute ranking scores by leveraging the full spectrum of trajectory information from multiple traversed paths ending at each candidate (as illustrated in Figure 2), our implementation currently utilizes only the most informative trajectory (i.e., the one with the longest traversed path) due to implementation complexity. Future work should explore adaptive methods to fully integrate the complete set of traversed paths into the candidate ranking process and compare the effectiveness of leveraging traversed paths at different levels.

References

Payal Chandak, Kexin Huang, and Marinka Zitnik. 2023. Building a knowledge graph to enable precision medicine. *Scientific Data*, 10(1):67.

Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.

Julien Delile, Srayanta Mukherjee, Anton Van Pamel, and Leonid Zhukov. 2024. Graph-based retriever captures the long tail of biomedical knowledge. *ArXiv*, abs/2402.12352.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *ArXiv*, abs/2404.16130.

Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. 2023a. Chatrec: Towards interactive and explainable llms-augmented recommender system. *arXiv preprint arXiv:2303.14524*.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023b. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.

Chunxi Guo, Zhiliang Tian, Jintao Tang, Shasha Li, Zhihua Wen, Kaixuan Wang, and Ting Wang. 2023. Retrieval-augmented gpt-3.5-based text-to-sql framework with sample-aware prompting and dynamic revision chain. In *International Conference on Neural Information Processing*, pages 341–356. Springer.

Haoyu Han, Yu Wang, Harry Shomer, Kai Guo, Jiayuan Ding, Yongjia Lei, Mahantesh Halappanavar, Ryan A Rossi, Subhabrata Mukherjee, Xianfeng Tang, et al. 2024. Retrieval-augmented generation with graphs (graphrag). *arXiv preprint arXiv:2501.00309*.

Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2025. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *Advances in Neural Information Processing Systems*, 37:132876–132907.

Minghao Hu, Yuxing Peng, Zhen Huang, and Dongsheng Li. 2019. Retrieve, read, rerank: Towards end-to-end multi-document reading comprehension. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. *Transactions on Machine Learning Research*.

Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji rong Wen. 2023. Structgpt: A general framework for large language model to reason over structured data. In *Conference on Empirical Methods in Natural Language Processing*.

- Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Suhang Wang, Yu Meng, and Jiawei Han. 2024a. Graph chain-of-thought: Augmenting large language models by reasoning on graphs. In *Annual Meeting of the Association for Computational Linguistics*.
- Bowen Jin, Yu Zhang, Sha Li, and Jiawei Han. 2024b. Bridging text data and graph data: Towards semantics and structure-aware knowledge discovery. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pages 1122–1125.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781.
- Oleksandr Kolomiyets and Marie-Francine Moens. 2011. A survey on question answering technology from an information retrieval perspective. *Information Sciences*, 181(24):5412–5434.
- Meng-Chieh Lee, Qi Zhu, Costas Mavromatis, Zhen Han, Soji Adeshina, Vassilis N. Ioannidis, Huzefa Rangwala, and Christos Faloutsos. 2024. Hybgrag: Hybrid retrieval-augmented generation on textual and relational knowledge bases. *ArXiv*, abs/2412.16311.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Millicent Li, Tongfei Chen, Benjamin Van Durme, and Patrick Xia. 2024. Multi-field adaptive retrieval. *arXiv preprint arXiv:2410.20056*.
- Lihui Liu, Yuzhong Chen, Mahashweta Das, Hao Yang, and Hanghang Tong. 2023. Knowledge graph question answering with ambiguous query. In *Proceedings of the ACM Web Conference 2023*.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*.
- Linhao Luo, Yuan-Fang Li, Reza Haf, and Shirui Pan. Reasoning on graphs: Faithful and interpretable large language model reasoning. In *The Twelfth International Conference on Learning Representations*.
- Yuning Mao, Pengcheng He, Xiaodong Liu, Yelong Shen, Jianfeng Gao, Jiawei Han, and Weizhu Chen. 2020. Generation-augmented retrieval for open-domain question answering. *arXiv preprint arXiv:2009.08553*.
- Yuning Mao, Pengcheng He, Xiaodong Liu, Yelong Shen, Jianfeng Gao, Jiawei Han, and Weizhu Chen. 2021. Generation-augmented retrieval for open-domain question answering. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*.
- Costas Mavromatis and George Karypis. 2024. Gnnrag: Graph neural retrieval for large language model reasoning. *ArXiv*, abs/2405.20139.
- Dang Nguyen, Viet Dac Lai, Seunghyun Yoon, Ryan A Rossi, Handong Zhao, Ruiyi Zhang, Puneet Mathur, Nedim Lipka, Yu Wang, Trung Bui, et al. 2024. Dynasaur: Large language agents beyond predefined actions. *arXiv preprint arXiv:2411.01747*.
- Bo Ni, Zheyuan Liu, Leyao Wang, Yongjia Lei, Yuying Zhao, Xueqi Cheng, Qingkai Zeng, Luna Dong, Yinglong Xia, Krishnaram Kenthapadi, et al. 2025. Towards trustworthy retrieval augmented generation for large language models: A survey. *arXiv preprint arXiv:2502.06872*.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Junhong Shen, Neil Tenenholtz, James Brian Hall, David Alvarez-Melis, and Nicolo Fusi. Tag-llm: Repurposing general-purpose llms for specialized domains. In *Forty-first International Conference on Machine Learning*.
- Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Sai Wang, Chen Lin, Yeyun Gong, Lionel M. Ni, Heung yeung Shum, and Jian Guo. 2023. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph. In *International Conference on Learning Representations*.
- Dhaval Taunk, Lakshya Khanna, Pavan Kandru, Vasudeva Varma, Charu Sharma, and Makarand Tapaswi. 2023. Grapeqa: Graph augmentation and pruning to enhance question-answering. *Companion Proceedings of the ACM Web Conference 2023*.
- Yijun Tian, Huan Song, Zichen Wang, Haozhu Wang, Ziqing Hu, Fang Wang, Nitesh V Chawla, and Panpan Xu. 2024. Graph neural prompting with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- MK Vijaymeena and K Kavitha. 2016. A survey on similarity measures in text mining. *Machine Learning and Applications: An International Journal*.
- Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. 2024. Knowledge graph prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19206–19214.
- Shirley Wu, Shiyu Zhao, Qian Huang, Kexin Huang, Michihiro Yasunaga, Kaidi Cao, Vassilis Ioannidis, Karthik Subbian, Jure Leskovec, and James Y Zou.

2025. Avatar: Optimizing llm agents for tool usage via contrastive reasoning. *Advances in Neural Information Processing Systems*, 37:25981–26010.

Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin Huang, Kaidi Cao, Qian Huang, Vassilis N Ioannidis, Karthik Subbian, James Zou, and Jure Leskovec. Stark: Benchmarking llm retrieval on textual and relational knowledge bases. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Yu Xia, Junda Wu, Sungchul Kim, Tong Yu, Ryan A Rossi, Haoliang Wang, and Julian McAuley. 2024. Knowledge-aware query expansion with large language models for textual and relational retrieval. *arXiv preprint arXiv:2410.13765*.

Mayi Xu, Yongqi Li, Ke Sun, and Tiejun Qian. 2024. Adaption-of-thought: Learning question difficulty improves large language models for reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5468–5495, Miami, Florida, USA. Association for Computational Linguistics.

Shangzi Xue, Zhenya Huang, Xin Lin, Jiayu Liu, Longhu Qin, Tianhuang Su, Haifeng Liu, and Qi Liu. 2024. Enhancing the completeness of rationales for multi-step question answering. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 2753–2763.

Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. 2019. End-to-end open-domain question answering with bertserini. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 72–77.

Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, and Li Yuan. 2023. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*.

Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. Qa-gnn: Reasoning with language models and knowledge graphs for question answering. In *North American Chapter of the Association for Computational Linguistics*.

Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *International conference on machine learning*. PMLR.

Huimin Zeng, Zhenrui Yue, Qian Jiang, and Dong Wang. 2024. Federated recommendation via hybrid retrieval augmented generation. *2024 IEEE International Conference on Big Data (BigData)*.

Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D Manning, and Jure Leskovec. Greaselm: Graph reasoning enhanced language models. In *International Conference on Learning Representations*.

Yanqiao Zhu, Yuanqi Du, Yinkai Wang, Yichen Xu, Jieyu Zhang, Qiang Liu, and Shu Wu. A survey on deep graph generation: Methods and applications. In *Learning on Graphs Conference*. PMLR.

A Summary of Notations

Table 6: Notations and the corresponding descriptions.

Notations	Definitions or Descriptions
\mathcal{B}	Text-rich Graph Knowledge Base (TG-KB)
$\mathcal{V}, \mathcal{E}, \mathcal{D}$	Set of Nodes, Categories and Documents of TG-KB
$\mathcal{D}_v, \mathcal{E}_v$	Document and Category of Node v
$Q \in \mathcal{Q}$	Query Q from Query set \mathcal{Q}
$Q^{\text{Struct}}, Q^{\text{Text}}$	Query targeted by structural and textual retrieval
$G = \{\mathcal{P}_i\}_{i=1}^{ \mathcal{G} }$	Planning Graph consisting of multiple reasoning paths
$\mathcal{P}_i = (p_{i1} \rightarrow \dots \rightarrow p_{iL_i})$	Reasoning path consisting of L_i sequential entities
$\mathcal{E}_{p_{ij}}, \mathcal{T}_{p_{ij}}$	Textual category and restriction of path entity p_{ij}
\tilde{C}	Retrieved candidates after reasoning module.
$\tilde{C}_i^l = \tilde{C}_i^{\text{Struct}} \cup \tilde{C}_i^{\text{Text}}$	Retrieved candidates at l^{th} layer for i^{th} path including structurally retrieved ones and textually retrieved ones.
C	Final retrieved candidates after organizing module.
$P_{\mathcal{Q} \times \mathcal{G}}$	Joint distribution of query and planning graph.
\mathcal{N}_v	Neighborhood of entity v
$\mathcal{I}_{p_{ij}}$	Traversal Identifier of Structural and Textual Retrieval
P_{Θ_1}	Planning module with its parameters Θ_1
P_{Θ_2}	Reasoning module with its parameters Θ_2
P_{Θ_3}	Organizing module with its parameters Θ_3

Dataset	# Entities	# Text Tokens	# Relations	Avg. Degree
AMAZON	1,035,542	592,067,882	9,443,802	18.2
MAG	1,872,968	212,602,571	39,802,116	43.5
PRIME	129,375	31,844,769	8,100,498	125.2

Table 7: Statistics of text-rich graph knowledge bases in STaRK benchmark (Wu et al.).

B Experimental Details

B.1 Datasets

To evaluate the effectiveness of our proposed framework, we conduct experiments using three Text-rich Graph Knowledge Bases (TG-KBs) from STaRK (Wu et al.). These TG-KBs cover a wide range of domains, including product reviews (Amazon), academic papers (MAG), and biomedical knowledge (Prime). Each TG-KB comprises a textual graph and an associated corpus, with the corpus containing documents linked to the nodes in the graph. Queries are meticulously crafted for each TG-KB and encompass varying levels of complexity, which desire different levels of textual and structural knowledge to answer.

Amazon: a dataset provides a realistic simulation of product search and recommendation. Its textual graph consists of four categories of nodes: *product*, *category*, *color*, and *brand*. Nodes are interconnected through relations such as *has_brand* and *has_category*. Textual documents encapsulate

Parameter	Value	Description
per_device_train_batch_size	4	Number of training samples per device.
gradient_accumulation_steps	8	Accumulates gradients over 8 steps.
num_train_epochs	100	Number of full passes on training dataset.
max_steps	1000	Maximum number of training steps; training stops once reached.
learning_rate	2e-4	Initial learning rate for the optimizer.
warmup_steps	5	Number of warm-up steps.
optim	adamw_8bit	8-bit AdamW for efficient training.
lr_scheduler_type	linear	Linear decay of learning rate.
weight_decay	0.01	L2 regularization to prevent overfitting.
seed	3407	Random seed for reproducibility.

Table 8: Hyperparameter for planning graph generator. properties of corresponding nodes, such as product descriptions and customer reviews.

MAG: a comprehensive resource for academic paper retrieval. In the textual graph, *papers* can be connected to other nodes, such as *field_of_study* via the *paper_has_topic_field_of_study* relation and *institution* through a combination of relations like *author_affiliated_with_institution* and *author_writes_paper*. Each *paper* document includes the title, abstract, and metadata, such as the publication date and venue, providing rich contextual knowledge for retrieval and analysis.

Prime: a highly domain-specific dataset. It focuses on medical inquiries and is sourced from the PrimeKG knowledge graph (Chandak et al., 2023), which comprises ten entity types and eighteen relation types, offering multiple target node categories, such as *disease*, *gene/protein*, and *drug*. The associated documents are aggregated from various databases, providing a rich and diverse source of medical knowledge.

Detailed dataset statistics are in Table 7.

B.2 Implementation Details

Planning Graph Generation: In Section 3.1, we follow previous works (Luo et al.; Wu et al.) to linearize the planning process by decomposing the planning graph into sequential reasoning paths, which can be generated by LLMs via next token prediction. Given the lack of ground-truth planning graphs for training, we prompt LLMs to synthesize these ground-truth planning graphs due to their superior reasoning capability. Since the nodes in planning graphs are entity categories/types, we include **Entity Type List** in the prompt. We also leverage in-context learning to help LLMs learn the mapping between the query and the corresponding planning graph. The detailed prompt for generating ground-truth planning graphs and parameters for fine-tuning the planning graph generator (i.e., Llama 3.2-3B) are shown in Prompt 1 and Table 8.

Trajectory Collection: As mentioned in Section 3.3, our reranker reorders the intermediate retrieved candidates based on their trajectory. To

Prompt 1: Planning Graph Generation

System Message: You are a planning graph finder agent. Your role is to:

1. Identify the underlying **meta-path** from a given question, which consists of the **entity types** at each reasoning step.
 2. Extract the **content restriction** for each **entity type** based on the question. If there is no restriction for an entity type, leave its value empty.
- You will be provided with a predefined **Entity Type List**. Only use the entity types from this list when constructing the meta-path and restrictions. Your response must be concise and strictly adhere to the specified **output format**.

Entity Type List: Provide the entity type list.

Demonstrations: Examples for in-context learning.

Output Format: *Metapath: ""*, *Restriction: {}*.

achieve this, we collect three key features: **Textual Fingerprint (TF)**, **Structural Fingerprint (SF)**, and **Traversal Identifier (TI)**.

Textual Fingerprint (TF): We record the BM25 similarity scores between the query and the traversed nodes computed. Since empirical observations indicate that the length of reasoning paths is typically less than three, we fix the textual fingerprint to the length of three by padding additional 0 similarity scores for those reasoning paths whose length is less than three, allowing for batch-wise training. Additionally, we append the initial semantic ranking score of the candidate computed using cosine similarity coupled with Ada-002 embedding to the end of three BM25-based similarity scores to complement the lexical perspective. This vector is then passed through a linear layer to be transformed into an embedding of size 256. Note that this initial ranking score is also used to select the intermediate retrieved candidates used for reranking.

Structural Fingerprint (SF): We concatenate the categories of all nodes in the corresponding reasoning path as a text sequence. If the reasoning path is shorter than three nodes, we prepend the sequence with "padding" tokens to ensure a fixed length. The structural fingerprint is then processed using a transformer model, which converts the sequence into an embedding of size 768, followed by a linear layer that projects it down to size 256.

Traversal Identifier (TI): We track whether each node is retrieved via textual matching or structural traversal and encoding them with distinct values by initializing a learnable embedding matrix mapping each traversal identifier encoding to a 3x256-dimensional embedding vector.

After obtaining all above three trajectory fea-

tures, we concatenate their obtained vectors into a unified vector ($256 + 256 + 256 \times 3 = 1280$) and apply two fully connected layers to transform the combined representation into a reranking score. This score determines the final ranking.

C Additional Results

C.1 Planning Performance

To investigate the performance of our planning graph generator, we conduct an experiment to evaluate the quality of the generated planning graphs for queries. We view the planning graph as correct if it enables the retrieval of at least one ground-truth candidate by traversal following its structure. In Table 9, most planning graphs are correct on MAG and Amazon datasets, demonstrating their effectiveness in guiding the traversal. Even for the domain-specific Prime dataset, more than 60% of planning graphs lead to correct retrieval, showing that LLMs still generate accurate plans and effectively support downstream reasoning and reranking.

Dataset	MAG	Amazon	Prime
Accuracy (%)	88.85	70.03	60.44

Table 9: Performance of planning graph generator on three datasets.

C.2 Query Pattern Analysis

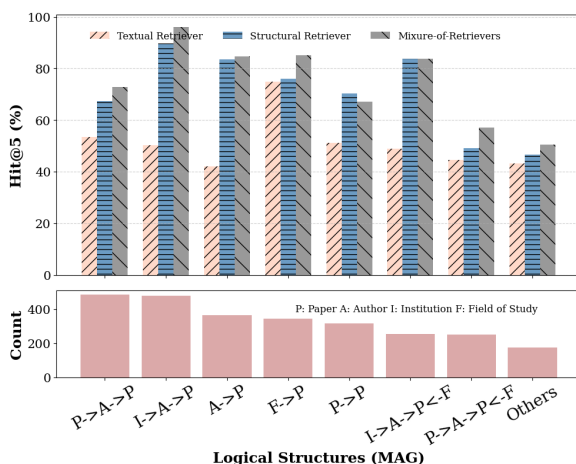


Figure 6: Imbalance number of queries and performance of different retrievers across different logic patterns.

Figure 6 illustrates the analysis of query patterns in the MAG dataset. With richer relational information, queries in this dataset form a wider variety of

patterns, including longer and more diverse structures. Similar to the Amazon dataset, we observe a general trend where the performance of MoR declines as the query count decreases across different patterns. Beyond this overall trend, certain query patterns in the MAG dataset stand out, such as "P → A → P" (Product-to-Author-to-Product) and "P → P" (Paper-to-Paper). Despite their relatively high occurrence, MoR still performs worse on these patterns. This is similar to low performance on the "Product → Product" pattern observed in the Amazon dataset, where repeated entity types appear within a single query. Such repetition causes the textual retriever to shift focus from the target to the repeated entities, leading to lower performance.

D Efficiency Analysis

D.1 Computational Complexity of MoR

Since our Planning-Reasoning-Organizing framework is multi-staged, we theoretically/empirically analyze the time complexity of each component:

- Planning:** the planning component generates the textual planning graph by linearized token generation using a pre-trained language model, which is the well-established technique for most LLMs. Since our planning graph is category-based and its textual description consists of simple terms extracted directly from the query, the resulting textual graph is relatively simple and maintains a short token length. Therefore, the generation is highly efficient and the time complexity is linear in the number of tokens, i.e., $\mathcal{O}(K \cdot D)$ with K as the number of planning graph tokens and D as the model dimension, following GPT-style decoding and key-value caching.
- Reasoning:** the reasoning component conducts layer-wise breadth-first traversal following generated textual planning graphs, the time complexity of which grows exponentially with the layers/hops of traversal. Let L be the maximum number of layers/hops across all reasoning paths, c be the total number of distinct reasoning paths decomposed from each planning graph, d be the average node degree, and t denote the number of nodes retrieved via textual similarity at each layer. We begin our mixed traversal from a set of s seed nodes and perform layerwise expansion by retrieving neighbors of nodes from the previous layer (constrained by the prescribed node

category at that corresponding layer in the planning graph). From the second layer onward, we additionally incorporate the top- t nodes ranked by textual similarity to the original query and these textually retrieved nodes also contribute neighbor expansion in the next layer. Assuming z represents the time used for textual matching and node entity constraint checking, the total time complexity for reasoning by following c reasoning paths is

$$\begin{aligned} & \mathcal{O}(c(\underbrace{(((sd+t)d+t)\dots)d+t}_{L})z) \\ &= \mathcal{O}(csd^{L-1}z + ct\frac{d^{L-1}-1}{d-1}z) \quad (8) \\ &\approx \mathcal{O}(Ad^{L-1}) \end{aligned}$$

as s, t, c, z are typically constants with small values compared with the exponentially growing paths and are absorbed together into A .

- **Organizing:** For each candidate given by the retrieved path from reasoning stage, the organization module calculates the ranking score, which takes $\mathcal{O}(Bd^{L-1})$ with d^{L-1} denotes the total number of reasoning paths and B denotes the time for calculating ranking score with the defined neural network for each path.

Taking the above analysis together, the time complexity of the entire framework is $\mathcal{O}(KD + (A + B)d^{L-1}) \approx \mathcal{O}(d^{L-1})$. Since our planning graph generation follows the typical LLM decoding mechanism, it does not introduce additional time overhead. Therefore, we do not focus on optimizing this stage. Instead, we optimize the Reasoning and Organizing stages, considering the exponentially growing number of retrieved paths:

- **Reasoning:** we parallelize the traversal process, as each reasoning path is explored independently.
- **Organizing:** we simultaneously evaluate a batch of reasoning paths and their candidates.

We report the theoretically analyzed and empirically verified time complexity for one query retrieval in Table 4. To account for differences across datasets and reduce the impact of stochasticity in the experiments, we perform three runs per dataset and report the average efficiency. We can see that despite the exponentially growing time complexity of reasoning and organizing, the empirical time consumption is rather low.

Dataset	Sequential (s)	Parallel (s)
Amazon	0.416	0.271
MAG	0.891	0.846
Prime	0.081	0.076

Table 10: Sequential vs. parallel processing time of different datasets.

D.2 Scalability of Mixed Traversal to Large-scale TG-KBs

Our reasoning process conducts layer-wise breadth-first traversal following generated textual planning graphs, the time complexity of which grows exponentially with the increase of traversal depth. The above analysis shows that the time cost (i.e., $\mathcal{O}(Ad^{L-1})$) is mainly governed by the average node degree d and the reasoning hop L . Therefore, even for large-scale TG-KBs, the traversal remains efficient since most graphs are sparse and the reasoning paths are reasonably short, both of which have been empirically verified in (Wu et al.).

Moreover, the three TG-KBs in our experiments are already large in scale while our reasoning traversal method achieves significant efficiency in addressing each query (see Table 7 and Table 10). Furthermore, thanks to the highly parallelizable nature of our traversal process, this runtime can be further reduced by leveraging multi-processing over traversing multiple reasoning paths.

E Extending to Question Answering Task

Prior studies have shown that improved retrieval often leads to improved question-answering outcomes (Mao et al., 2020; Lewis et al., 2020). To verify this positive correlation, we additionally evaluate the QA performance based on the retrieved information. Specifically, we use LLM GPT-4o as a judge to assess which retrieval method provides more helpful context for answering the question, allowing us to examine whether improvements in retrieval quality can convert to better QA performance. To mitigate hallucination, we enable LLM to respond with "I don't know" when neither retrieved candidate supports an answer. We sample 100 queries for each of the three TG-KBs and retrieve top-5 candidates by using both our proposed MoR and the classical textual matching baseline BM25. To mitigate the impact of positional bias, we evaluate two different presenting orders for candidates retrieved by the two methods, resulting in 200 evaluations per dataset. The win ratios (i.e., the percentage of preferences) for both methods across three datasets are reported in Table 11.

Dataset	BM25	MoR
MAG	27.0%	64.0%
Prime	33.0%	56.5%
Amazon	40.0%	42.5%

Table 11: Comparison of BM25 and MoR QA performance across datasets.

F Comprehensive Related Work

F.1 Retrieval-augmented Generation (RAG)

With the unprecedented success of recent LLMs in approaching human-level intelligence, retrieving relevant knowledge to support downstream generation has become increasingly crucial. Retrieval-augmented generation enhances generative tasks by integrating relevant information from external knowledge sources (He et al., 2025; Gao et al., 2023b; Han et al., 2024) and has been widely adopted to improve question-answering (Liu et al., 2023). In the context of LLMs, RAG has been utilized to mitigate hallucinations (Yao et al., 2023), enhance interpretability (Gao et al., 2023a), and enable dynamic knowledge updates (Wang et al., 2024). This work leverages RAG to retrieve supporting entities from TG-KBs, providing contextual grounding for answer generation. Depending on the type of knowledge retrieved, existing retrievers can be classified into structural and textual retrieval approaches, which are reviewed next.

F.2 Textual and Structural Retrieval

Since real-world knowledge is commonly stored in both textual and structural formats (Kolomiyets and Moens, 2011), such as indexed texts and knowledge graphs, each requires a retrieval method tailored to its unique representation. Textual retriever retrieves knowledge based on its similarity to the given query and can be categorized into: lexical methods (e.g., TF-IDF and BM25 (Robertson et al., 2009)) and semantic methods (e.g., DPR and Contriever (Karpukhin et al., 2020; Izacard et al., 2022)). Despite their broad applicability, the predefined linguistic rules and embedding-based semantics may struggle to capture the structural knowledge stored in graph-structured knowledge bases such as knowledge graphs and text-rich networks. To address this challenge, structural retrieval has been proposed by using graph analysis techniques (e.g., graph traversal (Wang et al., 2024; Jiang et al., 2023; Zhang et al.; Edge et al., 2024))

and graph machine learning models (e.g., graph neural networks (Yasunaga et al., 2021; Mavroumatis and Karypis, 2024)). Early methods extract local subgraphs of seeding nodes (Yasunaga et al., 2021; Taunk et al., 2023) or pre-define paths approaching answers (e.g., shortest paths (Luo et al.; Delile et al., 2024)). To avoid exponentially expanding neighbors in the local subgraphs and break the rigid logic routined by pre-defined paths, recent advancements integrated LLMs to dynamically adjust graph traversal (Sun et al., 2023; Wang et al., 2024; Jin et al., 2024a). While promising, frequently invoking LLMs introduces prohibitive resource overhead. Despite the above advancements in both textual and structural retrieval, they are often applied independently and fail to mutually reinforce each other. This motivates the recent research trend of developing hybrid retrieval, which is reviewed next.

F.3 Hybrid Retrieval

Recently, several works have explored hybrid knowledge retrieval from TG-KBs. One approach (Xia et al., 2024; Li et al., 2024) aggregate documents from neighboring nodes, with Xia et al. (2024) applying relational filtering to remove irrelevant neighbors and Li et al. (2024) weighting neighbors based on field importance. Another approach (Lee et al., 2024) uses LLMs to choose either structural or textual retrieval. In contrast, our proposed MoR fully leverages the graph structure and rich texts by integrating textual matching and graph traversal into a unified framework, enabling a more seamless and interpretable interaction between structural and textual knowledge