

# RiTeK: A Dataset for Large Language Models Complex Reasoning over Textual Knowledge Graphs in Medicine

Anonymous ACL submission

## Abstract

Answering complex real-world questions in the medical domain often requires accurate retrieval from medical Textual Knowledge Graphs (medical TKGs), as the relational path information from TKGs could enhance the inference ability of Large Language Models (LLMs). However, the main bottlenecks lie in the scarcity of existing medical TKGs, the limited expressiveness of their topological structures, and the lack of comprehensive evaluations of current retrievers for medical TKGs. To address these challenges, we first develop a Dataset<sup>1</sup> for LLMs Complex Reasoning over medical Textual Knowledge Graphs (RiTeK), covering a broad range of topological structures. Specifically, we synthesize realistic user queries integrating diverse topological structures, relational information, and complex textual descriptions. We conduct a rigorous medical expert evaluation process to assess and validate the quality of our synthesized queries. RiTeK also serves as a comprehensive benchmark dataset for evaluating the capabilities of retrieval systems built upon LLMs. By assessing 11 representative retrievers on this benchmark, we observe that existing methods struggle to perform well, revealing notable limitations in current LLM-driven retrieval approaches. These findings highlight the pressing need for more effective retrieval systems tailored for semi-structured data in the medical domain.

## 1 Introduction

Although large language models (LLMs) have made significant strides in natural language processing (NLP), complex question answering still remains a challenge. Medical professionals, for instance, often need to express complex information that combines flexible inputs with specific,

<sup>1</sup>The dataset is available here: [https://anonymous.4open.science/r/Riteck\\_submission\\_version-026B/readme.md](https://anonymous.4open.science/r/Riteck_submission_version-026B/readme.md)

structured constraints. Consider the query, “Which organ or tissue function that circulates maternal and fetal blood is affected by Fetal Distress?” compared with the simpler version, “What does Fetal Distress affect?” Accurately addressing such complex queries is crucial, as it directly impacts healthcare diagnosis and treatment planning.

To effectively answer these queries, organizing the underlying knowledge using medical TKGs becomes essential. TKGs integrate unstructured data, such as textual descriptions of nodes (e.g., the definition of the medical term *Placental Circulation*) with structured data, like the relationships between entities within the graph (e.g., the relationship between *Fetal Distress* and *Placental Circulation* is *affects*). This integration enables TKGs to represent comprehensive knowledge tailored to specific applications, rendering them invaluable, especially in the medical field, where accuracy and reliability are critically important.

However, existing datasets (Wu et al., 2024b,a) exhibit several critical limitations: they are overly simplistic, typically limited to 1-2 hop reasoning paths; they lack diverse topological structure templates<sup>2</sup> and rich relation types; or they fail to incorporate complex constraints<sup>3</sup>. Consequently, these datasets are inadequate for capturing the complexity of retrieval tasks involving medical TKGs, where queries demand multi-hop reasoning, diverse topological structure templates, and multiple interdependent constraints. Moreover, the absence of textual properties in existing medical TKGs limits their effectiveness in delivering comprehensive answers.

To bridge this gap, we introduce RiTeK, a large-

<sup>2</sup>The details of the topological structure are provided in the Appendix C.

<sup>3</sup>Constraints are particularly important in KBQA as they help filter out irrelevant information from large knowledge bases, narrowing the search space and improving both efficiency and accuracy.

076 scale dataset for complex reasoning over medical  
077 TKGs. In this progress, one primary technical  
078 challenge we address is the accurate simulation  
079 of user queries with different reasoning types (e.g.,  
080 six topological structures in Figure 1 within medi-  
081 cal TKGs, ensuring that these queries are relevant  
082 and reflective of real-world medical scenarios in-  
083 volving patients, doctors, and medical scientists.  
084 This challenge stems from the interdependence be-  
085 tween textual and relational information, the in-  
086 herent complexity of medical terminology and re-  
087 lationships, and the limited availability of textual  
088 descriptions for medical terms. We refer to the  
089 framework of Wu et al. (2024b) to simulate user  
090 queries and construct precise ground-truth answers.  
091 We incorporate richer topological structures that  
092 extend beyond the traditional 2- and 3-hop reason-  
093 ing patterns to better reflect real-world medical  
094 scenarios. Compared with datasets like BioKG-  
095 Bench (Lin et al., 2024) and PrimeKGQA (Yan  
096 et al., 2024), the queries in RiTek not only con-  
097 tains structural information that requires reasoning  
098 but also includes textual information related to the  
099 ground-truth answer, making the task more chal-  
100 lenging. In addition, we enrich the textual descrip-  
101 tions of each node and incorporate more ontology  
102 structure.

103 The key features of RiTeK are summarized as  
104 follows: (1) it integrates rich ontological structures  
105 and comprehensive textual descriptions, with con-  
106 tent quality rigorously validated by medical experts  
107 to ensure reliability; (2) it constructs queries that  
108 capture complex relational dependencies and nu-  
109 anced linguistic variations; and (3) it introduces  
110 context-sensitive reasoning, where effective re-  
111 trieval depends not only on a model’s reasoning  
112 ability but also on its semantic alignment with en-  
113 tity constraints embedded in the query.

114 Moreover, we systematically investigate the per-  
115 formance of existing retrieval systems on RiTeK  
116 and provide insights to guide future research. In  
117 particular, we identify key challenges in process-  
118 ing textual and relational data with complex on-  
119 tology structures and in mitigating latency issues  
120 on large-scale SKBs containing millions of entities  
121 and relations.

## 122 2 Related Work

### 123 Datasets of Question Answering over Document.

124 This area of research centers on extracting answers  
125 from document sources (Rajpurkar, 2016; Dunn

126 et al., 2017; Joshi et al., 2017; Trischler et al., 2016;  
127 Welbl et al., 2018; Yang et al., 2018; Jin et al.,  
128 2021, 2019; Hendrycks et al., 2020). For exam-  
129 ple, SQuAD (Rajpurkar, 2016) assesses a model’s  
130 ability to interpret and retrieve answers from a sin-  
131 gle document, focusing on comprehension within  
132 a defined context. PubMedQA (Jin et al., 2019) tar-  
133 gets reasoning over complex biomedical literature.  
134 MedQA-CS (Yao et al., 2024b) aims to simulate  
135 authentic medical examination scenarios in clini-  
136 cal education. However, existing unstructured QA  
137 datasets often lack the depth required for relational  
138 reasoning and failed to address complex user in-  
139 quiries. In contrast, our research involves queries  
140 that demand more complex relational reasoning,  
141 challenging the model’s ability to navigate and uti-  
142 lize structured information effectively.

### 143 Datasets of Question Answering over Knowl- 144 edge Graph.

145 Structured QA datasets challenge  
146 models to retrieve answers from knowledge graphs,  
147 which serve as structured databases for factual rea-  
148 soning (Zhang et al., 2018; Yih et al., 2016; Gu  
149 et al., 2021; Bao et al., 2016; Trivedi et al., 2017).  
150 For instance, MetaQA (Zhang et al., 2018) requires  
151 models to infer multi-hop relational paths across  
152 entities. To test the models’ abilities to decompose  
153 the constraint information in the queries, WebQues-  
154 tionsSP (Yih et al., 2016) is proposed. GraiQA (Gu  
155 et al., 2021) aims to facilitate the answering of more  
156 complex questions, as it allows queries to involve  
157 up to four relations and optionally includes func-  
158 tions such as counting, superlatives, and compara-  
159 tives. However, these datasets primarily focus on  
160 relational information; the absence of textual con-  
161 text restricts query diversity and limits the semantic  
162 expressiveness of reasoning within predefined re-  
163 lationships and entities.

### 164 Datasets of Question Answering over Textual 165 Knowledge Graph.

166 To integrate textual informa-  
167 tion into knowledge graphs and queries, the STaRK  
168 dataset (Prime, Amazon, Mag) (Wu et al., 2024b)  
169 was proposed. To the best of our knowledge,  
170 STaRK remains the only dataset that integrates re-  
171 lational and textual information for question ans-  
172 wering over TKGs. However, this dataset exhibits lim-  
173 ited topological structure coverage, which restricts  
174 its ability to handle complex multi-hop queries, par-  
175 ticularly in the medical domain. Furthermore, the  
176 lack of detailed node descriptions further impairs  
177 a model’s ability to comprehend query semantics.  
RiTeK addresses these limitations by incorporating  
richer topological structures and more extensive

textual information into both knowledge graphs and queries. This integration leads to more comprehensive and nuanced responses, providing deeper insights drawn from abundant textual data.

### 3 Problem Statement

**Textual Knowledge Graph** A Textual Knowledge Graph (TKG) is defined as a graph  $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{D})$ , where  $\mathcal{E}$  denotes a set of entities and  $\mathcal{R}$  denotes a set of relations among these entities. In a TKG, the entities and relations are usually organized as *facts*, and each fact is defined as a triplet  $(h, r, t)$ , where  $h, t \in \mathcal{E}$  and  $r \in \mathcal{R}$  denote the head entity, tail entity and the relation between the two entities, respectively. Each entity  $e$  ( $e = h$  or  $e = t$ ) in  $\mathcal{G}$  is associated with a textual document  $d^e \in \mathcal{D}$  describing its information.

**Complex Question Answering over Textual Knowledge Graph** Given a textual knowledge graph  $\mathcal{G}$  and input query  $q$ , the model is expected to generate the answers  $a \in \mathcal{E}$ , which satisfy the relational constraints defined by the structure of  $\mathcal{G}$  as specified in  $q$ , and the associated document  $d^e$  needs to satisfy the knowledge required to solve  $q$ .

**Textual Triple Graph** Unlike traditional knowledge graphs, where each node represents an entity and each edge denotes the relationship between nodes, in the textual triple graph, each node corresponds to a triple (head entity, relation, tail entity) along with the textual description of each entity. In this context, the relation indicates whether the two triples are connected. To be specific, let  $\mathcal{G}^* = (V, E)$  denote a graph consisting of a set of nodes  $V$  and a set of edges  $E \in V \times V$ . We denote by  $n$  the number of nodes in  $\mathcal{G}$  and by  $m$  its number of edges. Each node  $v = (h, r, t, T(h), T(t)) \in V$ , where  $T(*)$  denotes the textual description of an entity.

## 4 Dataset for LLMs Complex Reasoning over TKGs (RiTeK)

### 4.1 Medical Textual Knowledge Graph Construction

We construct two medical TKGs based on PharmKG (Zheng et al., 2021) and ADint (Xiao et al., 2024), as the increased number of entity and relation types introduces significant challenges for path retrieval in question answering over textual knowledge graphs. To enrich entity representations, we incorporate textual details from various

TKG Dataset	# Entities	# Relation	# Triple	# Coverage
STaRK-Amazon	4	4	9,443,802	–
STaRK-Mag	4	4	39,802,116	–
STaRK-Prime	10	18	8,100,498	15.29%
RiTeK-PharmKG	3	29	500,958	95.61%
RiTeK-ADint	102	15	1,017,284	36.73%

Table 1: Dataset Statistics of constructed medical textual knowledge graphs. # Coverage refers to the proportion of nodes with textual descriptions. # Entities denotes the number of entity types, and # Relations indicates the number of relation types. As the textual information of the provided nodes is difficult to quantify statistically, we do not include the corresponding statistics for STaRK-Mag and STaRK-Amazon.

databases, including Ensembl, UMLS, and Mondo Disease Ontology. As shown in Table 1, our constructed TKGs provide greater node textual coverage, as well as a larger variety of entity and relation types. For further details on these two medical TKGs, please refer to Appendix A.4.

### 4.2 Question Answering Dataset Construction

QA Dataset	# queries	# topological structure	# instance rate	train/val/test
STaRK-Amazon	9,100	1	4	0.65/0.17/0.18
STaRK-Mag	13,323	4	1.25	0.60/0.20/0.20
STaRK-Prime	1,1204	3	9.3	0.55/0.20/0.25
RiTeK-PharmKG	1,0235	6	11.33	0.80/0.10/0.10
RiTeK-ADint	5322	6	9.67	0.80/0.10/0.10

Table 2: Statistical Overview of the Textual KBQA benchmark Datasets. Instance rate refers to the average number of relational templates per topological structure.

#### 4.2.1 Overview

We developed two question-answering datasets, **RiTeK-PharmKG** and **RiTeK-ADint**, based on textual knowledge graphs for complex reasoning. These datasets notably feature queries that integrate relational and textual knowledge, incorporating relational templates with broader coverage and higher instance rates. Additionally, to enhance their applicability in practical scenarios, these queries mimic real-world query patterns, exhibiting a natural-sounding quality and flexible formats. Specifically, RiTeK-PharmKG consists of 10,235 synthesized queries. To maximize the coverage of different question topologies, we generate the queries following the six types of topological structure (e.g., multi-hop and constrained multi-hop). For the synthesized queries, we developed 68 relational templates, crafted by medical experts and detailed in

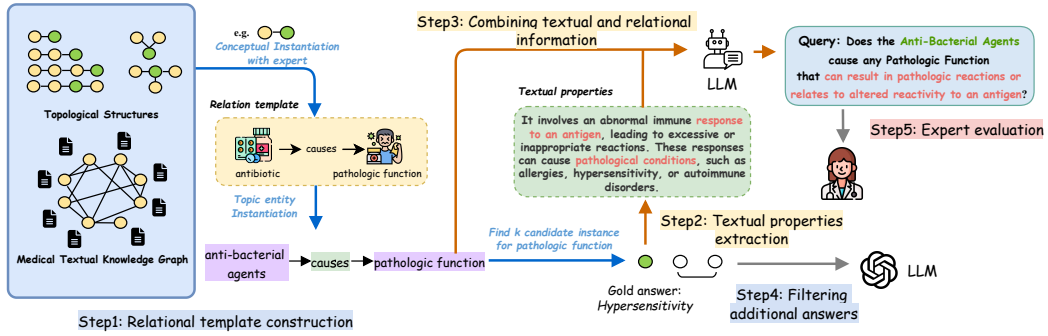


Figure 1: The process of constructing textual structured retrieval datasets involves five main steps, 1) Relational template construction: Create the relation template for TKG using the expert-designed topological structure. 2) Extract Textual Properties: Choose one node as the answer node that meets the relational requirement, and extract relevant textual properties. 3) Combine Information: Merge the relational information and textual properties to form a natural-sounding query. 4) Filtering additional answers: Check if the left nodes satisfy the textual properties to establish other ground truth nodes. 5) Expert Evaluation: The medical experts evaluate the naturalness, diversity, and practicality of the dataset.

Appendix B.1, to encompass various relation types and ensure practical relevance. The instance rate of 11.33, which is higher than that of the current TKG dataset STaRK (Amazon, Mag, and Prime), highlights the greater diversity of our dataset. RiTeK-ADint consists of 5322 synthesized queries and covers 6 topological structures, with 58 relational templates. Further details are provided in Appendix B.2. To capture the diverse language styles used by different users, we follow STaRK and simulate three distinct roles: medical scientist, doctor, and patient. We divide the synthesized queries on each dataset into training, validation, and test subsets, with the ratios detailed in Table 2. Further details on the scale of our QA benchmarks can be found in Table 2.

#### 4.2.2 Construction Pipeline

We present the pipeline used to generate large-scale medical QA datasets on TKGs. The core idea is to intertwine relational information and textual properties within the queries, accurately constructing ground-truth answers that exhibit more complex topological structures. The construction of the QA datasets (Figure 1) generally involves five steps, and the specific processes vary depending on the characteristics of each dataset. These steps are as follows.

**Relational Template Construction.** As shown in Figure 1 Step 1, we first created templates based on the 6 designed topological structures (Li and Ji, 2022), which were evaluated by medical experts to ensure their practical relevance and value. Afterward, the topological structures are instanti-

ated conceptually with experts. For instance, for the topological structure *Head entity–relation–tail entity*, the *"(antibiotic) causes <pathologic function>"* is a valid and common medical relation template, as antibiotics, particularly penicillin and cephalosporins, are well-known for triggering drug hypersensitivity reactions. This makes it a medically reasonable and frequently observed relationship. We then converted these relation templates into specific relationship queries, such as *"Anti-Bacterial Agents causes pathologic function."* Since each query could correspond to one or more candidate entities, we matched the queries with the textual KG to obtain k candidate entities.

**Extracting Textual Properties.** As shown in Figure 1 Step 2, for the k candidate answers that meet the relationship criteria, we select one entity as the *gold answer* and use GPT-4 to extract textual properties from the entity’s associated document. For instance, in the relationship *"Anti-Bacterial Agents causes pathologic function,"* we selected *"Hypersensitivity"* as the gold answer and extracted its textual properties. These textual properties elaborate on the concept of hypersensitivity, highlighting its key characteristics, which make it more likely to meet the inquirer’s needs.

**Combining Textual and Relational Information.** As shown in Figure 1 Step 3, after obtaining the relational templates and textual properties, we combine these components to synthesize the queries. We chose GPT-4 as the LLM for query synthesis, as it excels at generating natural, human-like questions. Additionally, we optimized the prompt and incorporated instructions for different personas to

make the queries more diverse and realistic. This approach enhances the quality of our dataset and increases the demands on our model’s reasoning capabilities. For details on using GPT-4 to generate this query, please refer to Appendix A.5.

**Filtering Additional Answers.** As shown in Figure 1 Step 4, in addition to the gold answer from which the textual properties are extracted, we need to evaluate whether other candidates meet the requirements of the query in order to include them in the final answer set. We use multiple LLMs to assess whether each candidate’s description satisfies the textual requirements of the query. Only candidates that pass validation by all LLMs will be added to the final answer set.

**Human Evaluation.** We invited four medical experts to evaluate 1,000 synthetic queries sampled from two datasets. The evaluation was conducted using a 5-point Likert-like scale across three dimensions. Naturalness measures how grammatically correct and human-like the queries sound. Diversity assesses whether the queries exhibit complex logical structures and encompass multiple entities, relations, and textual requirements. Practicality evaluates the real-world applicability of the generated queries and their likelihood of being encountered in real clinical or everyday scenarios.

The scores were ultimately converted into percentages representing the rates of Positive and Acceptable responses. We found that the evaluation results provided by GPT-4 for our generated dataset were largely consistent with assessments from medical experts. For shorter queries, such as “What gene is inhibited by naloxone?”, GPT-4 noted the limited relational and textual information contained within and consequently assigned a lower Diversity score. Both GPT-4 and medical experts agreed that certain rare relationship types, such as “an ancestor of”, are infrequently encountered in everyday Q&A scenarios and are more common in medical education contexts. Only a very small number of queries exhibited issues with insufficient Practicality. The results of this evaluation are summarized in Table 3. The data in the table represents the Positive/Acceptable rates (%) from GPT-4.

	Naturalness	Diversity	Practicality
RiTeK-PharmKG	81.80/99.60	81.6/99.40	67.4/97.8
RiTeK-ADint	81.20/99.20	74.80/100	68.60/96.60

Table 3: Positive/Acceptable rates(%) from experts

### 4.2.3 Data Distribution Analysis

We chose Shannon Entropy and Type-Token Ratio (TTR) as metrics to evaluate query diversity generated in our two datasets. Shannon Entropy takes into account the frequency of each word, measuring the evenness of word distribution in the text, while Type-Token Ratio reflects the variety of words, with a higher value indicating greater diversity in the generated queries. We found the TTR values for both RiTeK-PharmKG and RiTeK-ADint surpass those of STaRK-Prime, demonstrating that the queries generated in our datasets exhibit high complexity and diversity (The results are shown in Appendix D and Table 6). For Shannon Entropy, our results are comparable to STaRK-Prime. Since our RiTeK-ADint dataset involves a wide range of non-pharmacological interventions (NPIs), lifestyle modifications, and environmental factors, it introduces a richer variety of specialized terminology and concepts into the synthesized queries. This expanded vocabulary diversity results in significantly higher Shannon Entropy compared to other medical domain datasets. However, since our two datasets are derived from the medical domain, the frequent repetition of specialized medical terminology, as well as the more concentrated vocabulary compared to general-domain texts, results in slightly lower Shannon Entropy for our datasets than for the other two general-domain datasets. For further analysis about the distribution of query lengths and answer length, please refer Appendix A.6.

## 5 Experiments

### 5.1 Retrieval Models and Evaluation Metrics

We evaluated 11 representative retrieval models on our benchmark datasets under both zero-shot and few-shot learning settings. In addition to our benchmark dataset, we also evaluated the models on STaRK-Prime (Wu et al., 2024b), a textual question answering dataset with minimal ontological structure in its queries, including:

- GPT-4 (Achiam et al., 2023): We use GPT-4 with the instruction to generate the answers directly.
- Random Walk (Lovász, 1993): Starting from the topic entity, a random walk algorithm is applied to explore paths in the textual triple graph in the maximum depth  $d$ .
- MCTS (Chaslot, 2010): Starting from the topic entity, a Monte Carlo tree search algorithm is applied to explore paths in the textual triple graph

	Approach	RiTeK-PharmKG						RiTeK-ADint						STaRK-Prime					
		Exact Match			Rouge-1			Exact Match			Rouge-1			Exact Match			Rouge-1		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Zero-Shot	GPT-4	11.39	10.90	11.03	15.56	15.50	15.30	7.26	12.10	8.03	13.71	27.64	16.35	5.23	6.81	4.65	11.31	16.35	11.31
	+Random Walk (Lovász, 1993)	12.27	11.86	11.96	14.69	14.15	14.30	15.12	22.68	16.52	20.87	32.92	23.25	7.50	8.20	6.48	13.90	17.31	13.32
	+MCTS (Chaslot, 2010)	17.17	16.54	16.68	19.09	18.44	18.60	16.97	24.41	18.35	22.82	34.69	25.20	7.64	8.36	6.52	14.04	17.45	13.38
	+COT (Wei et al., 2022)	13.11	16.42	13.70	17.53	22.57	18.40	10.52	19.78	11.95	17.79	37.25	20.97	6.47	8.23	5.81	12.61	17.99	12.47
	+TOT (Yao et al., 2024a)	7.31	7.32	7.22	13.21	14.67	13.42	3.97	9.65	5.28	12.90	25.44	15.96	2.99	3.08	2.55	9.50	9.81	8.65
	+GOT (Besta et al., 2024)	3.56	4.20	3.75	10.86	11.84	11.06	2.61	3.32	2.81	15.09	17.63	15.84	1.99	2.20	1.78	9.89	9.34	8.72
	+TOG (Sun et al., 2023)	29.85	<b>38.19</b>	<b>31.14</b>	31.38	<b>40.37</b>	<b>32.92</b>	23.08	<b>40.63</b>	25.81	27.81	<b>48.93</b>	31.54	12.14	<b>15.76</b>	<b>11.27</b>	<b>18.67</b>	<b>24.75</b>	<b>18.42</b>
	+G-retriever (He et al., 2024)	11.21	13.39	11.60	15.01	18.54	15.62	10.97	19.05	12.52	17.27	32.99	20.41	6.23	6.61	5.17	12.01	14.92	11.40
	+KAR (Xia et al., 2024)	<b>30.95</b>	23.99	25.18	<b>33.65</b>	26.11	27.50	<b>39.59</b>	24.00	<b>27.29</b>	<b>46.54</b>	28.87	<b>32.80</b>	<b>12.02</b>	14.49	11.12	18.04	22.20	17.61
Few-Shot	GPT-4	13.75	15.54	14.04	16.84	19.84	17.49	17.57	17.91	17.48	25.50	28.08	26.04	7.79	6.41	5.91	14.03	13.53	12.14
	+Random Walk (Lovász, 1993)	11.02	13.28	11.32	14.46	17.88	14.92	22.99	22.79	22.75	29.10	29.07	28.95	9.93	6.93	7.34	16.54	13.02	13.45
	+MCTS (Chaslot, 2010)	17.79	17.11	17.30	20.97	20.29	20.48	19.51	27.32	20.91	24.71	36.25	26.96	9.57	6.89	7.14	15.92	12.55	12.88
	+COT (Wei et al., 2022)	17.29	16.91	16.99	21.55	20.97	21.13	18.57	18.12	18.26	26.68	26.62	26.53	8.13	5.91	5.99	14.03	13.53	12.14
	+TOT (Yao et al., 2024a)	14.74	14.74	14.63	19.22	19.14	18.97	13.28	13.17	13.21	24.65	24.72	24.60	12.84	10.11	10.36	6.93	4.85	5.06
	+GOT (Besta et al., 2024)	12.10	12.22	12.06	17.38	17.31	17.19	15.84	15.32	15.42	26.20	25.89	25.91	5.37	3.73	3.78	12.69	9.98	10.17
	+TOG (Sun et al., 2023)	<b>29.14</b>	<b>42.33</b>	<b>32.36</b>	<b>30.40</b>	<b>44.00</b>	<b>33.88</b>	<b>26.50</b>	<b>47.13</b>	<b>33.83</b>	<b>29.46</b>	<b>49.69</b>	<b>36.43</b>	<b>14.41</b>	<b>20.39</b>	<b>16.40</b>	<b>19.75</b>	<b>26.61</b>	<b>20.14</b>
	+G-retriever (He et al., 2024)	12.51	12.14	12.22	15.94	15.44	15.57	17.47	17.50	17.32	24.87	24.92	24.71	7.72	5.75	5.86	14.63	11.92	12.10
	+KAR (Xia et al., 2024)	27.35	27.43	26.99	29.74	29.76	29.34	34.68	33.42	33.48	40.15	38.55	38.88	13.01	15.50	12.21	19.00	23.10	18.00
Supervised	G-retriever (He et al., 2024)	38.71	37.11	37.62	39.78	39.18	39.31	47.93	47.16	47.41	54.68	54.00	54.24	16.14	16.47	14.11	17.21	27.86	19.21
	GCR (Luo et al., 2024)	44.38	<b>57.28</b>	47.71	46.04	<b>58.83</b>	49.44	43.52	<b>60.78</b>	48.07	49.47	<b>65.57</b>	54.24	<b>19.03</b>	<b>26.89</b>	<b>18.94</b>	<b>28.01</b>	<b>37.18</b>	<b>28.75</b>
	GNN-RAG (Mavromatis and Karypis, 2024)	<b>50.78</b>	49.28	<b>49.72</b>	<b>51.66</b>	50.29	<b>50.73</b>	<b>51.04</b>	50.59	<b>50.55</b>	<b>56.49</b>	56.09	56.09	16.00	15.04	14.50	24.78	23.51	22.99

Table 4: Results of various approaches for question answering with complex reasoning on RiTeK-PharmKG, RiTeK-ADint and STaRK-Prime. P refers to the Precision, R refers to the recall. In the experiments, the GPT-4 version is GPT4o-mini.

in the maximum depth  $d$ . In this work, we set the  $d = 3$ .

- Chain-of-Thought (COT) (Wei et al., 2022): We designed the instruction to guide GPT-4 in generating the answer step by step, with the output formatted as step-by-step reasoning: explanation, answer: medical terms.
- Tree-of-Thought (TOT) (Yao et al., 2024a): We structured the reasoning process as a tree search, where multiple intermediate reasoning paths are explored in parallel. GPT-4 evaluates and expands promising paths based on a voting or scoring mechanism.
- Graph-of-Thought (GOT) (Besta et al., 2024): We represented the reasoning process as a graph structure, where nodes capture different reasoning states and edges denote transitions. GPT-4 traverses the graph to aggregate information and synthesize the final answer.
- Think-on-Graph (TOG) (Sun et al., 2023): is a reasoning framework that enables large language models to interactively perform beam search over knowledge graphs, discovering and evaluating promising reasoning paths without additional training.
- G-retriever (He et al., 2024): A RAG-based approach that retrieves query-relevant subgraphs using the Prize-Collecting Steiner Tree (PCST) algorithm to enhance LLM understanding and reasoning over textual graphs.
- KAR (Xia et al., 2024): A knowledge-aware query expansion method that augments LLMs with structured document relations from a knowl-

edge graph, using relation-aware filtering to improve retrieval for semi-structured queries.

We evaluated the 3 representative retrieval models on our benchmark datasets and STaRK-Prime under supervised learning settings, including:

- G-retriever (He et al., 2024): A RAG-based approach that retrieves query-relevant subgraphs using the Prize-Collecting Steiner Tree (PCST) algorithm to enhance LLM understanding and reasoning over textual graphs.
- GCR (Luo et al., 2024): A knowledge-aware query expansion method that augments LLMs with document-based relational signals to improve retrieval for semi-structured queries.
- GNN-RAG (Mavromatis and Karypis, 2024). A method that uses a GNN to retrieve relevant answers and extract the shortest paths connecting the topic entity and answers, which are then verbalized and fed into the LLM to enhance retrieval-augmented generation (RAG) performance.

We evaluated the outputs of different methods using several metrics, including Exact Match (EM) (Rajpurkar, 2016; Li and Huang, 2023), which assesses whether the predicted sequence exactly matches the reference, awarding credit only for perfect matches. Additionally, we employed ROUGE-1 (Cohan and Goharian, 2016) to measure unigram overlap between the predicted and reference sequences, providing partial credit for shared words even when the sequences are not identical. To ensure fairness in the comparison, the instructions and examples are the same for both the zero-

shot and few-shot settings, respectively.

## 5.2 Results and Discussion

Table 4 shows the experiment results of various approaches based on Exact Match and Rouge-1. We have the following observations. Zero-shot and few-shot setting: (1) We observed that the baseline models struggle to generate the correct answers on RiTeK-PharmKG and RiTeK-ADint. For GPT-4 and GPT+COT, they are challenges in utilizing reasoning information from the graph. Although GPT+COT can utilize step-by-step reasoning, it still relies on the inherent knowledge of the LLM, which limits its ability to apply clear logical reasoning based on knowledge graphs. For the Random Walk, while it can provide reasoning paths, its random nature limits its ability to accurately identify the correct path information. However, it could get the better performance than GPT-4 in RiTeK-ADint and STaRK-Prime in the zero/few-shot setting. (2) Tree-of-Thought (ToT) and Graph-of-Thought (GoT) attempt to guide LLM reasoning through structured prompting, encouraging step-by-step or graph-based logical thinking. However, on complex textual KBQA datasets like RiTeK-PharmKG and RiTeK-ADint, both methods consistently underperform, with F1 scores far below those of retrieval-augmented approaches like KAR (e.g., ToT: 13.42 vs. KAR: 27.50 in zero-shot). This suggests that the internal knowledge and reasoning capabilities of LLMs alone are insufficient for tasks that require fine-grained relational understanding and the integration of attribute information from the query. Despite their logical scaffolding, ToT and GoT struggle to recover factual precision without access to external structured knowledge. (3) KAR achieves strong performance on medical datasets like RiTeK-PharmKG and RiTeK-ADint, outperforming baselines in both zero-shot and few-shot settings. Its main strength lies in combining textual semantics with structured KG relations to generate accurate and context-aware query expansions. However, KAR relies on retrieving the top-n relevant documents; however, determining an appropriate value for n and the optimal order in which to select documents is non-trivial. (4) G-Retriever shows moderate performance across medical datasets, but generally underperforms compared to methods like KAR or TOG in both zero-shot and few-shot settings. For example, on STaRK-Prime, its ROUGE-1 F1 score (5.17 vs. 11.12 zero-shot) lags significantly behind

KAR, This indicates a weaker ability to handle complex relational constraints, particularly when the answer’s attributes are embedded in the query. Its main strength lies in interpretable subgraph selection using PCST, which enhances explainability and helps mitigate hallucinations. (5)TOG performs moderately in zero-shot settings but shows strong gains in few-shot scenarios, achieving top-tier ROUGE-1 F1 scores like 37.11 on RiTeK-ADint and 36.43 on STaRK-Prime. This highlights its ability to leverage demonstrations to guide accurate reasoning over knowledge graphs, especially in complex biomedical tasks.

In the setting of supervised fine tuning, GCR achieves the best overall performance across all three medical benchmarks in the supervised setting, with scores like 57.28 ROUGE-1 F1 on ADint and 49.72 on STaRK-Prime, demonstrating its strength in generating faithful, KG-grounded answers. However, GCR relies on pre-constructed KG-Trie indices. We found that GNN-RAG achieves better performance on the RiTeK-PharmKG and RiTeK-ADint datasets, demonstrating its ability to retrieve relevant path information from the graph. However, since it primarily relies on shortest paths, it may overlook critical reasoning information embedded in more complex or indirect graph structures.

## 5.3 Analysis

### 5.3.1 Effect of Different LLMs on Retriever Effectiveness

In this part, we analyze the influence of different LLMs on the retrievers. Table 5 presents the performance of three retrieval settings, G-retriever, GNN-RAG, and without retriever, in three LLMs of the backbone: Llama 3.1 8b, Llama2-chat-7b, and Biomixtral 7b, on three datasets. Overall, G-retriever consistently outperforms other approaches across most metrics, particularly in Rouge-1 F1 scores. For instance, on RiTeK-ADint, G-retriever with Llama 3.1 8b achieves the highest F1 score of 56.87, while the GNN-RAG and no retriever baselines lag behind. Similarly, G-retriever reaches 55.02 F1 on Biomixtral for the same dataset, showcasing its robustness across model sizes. In contrast, GNN-RAG shows variable performance, sometimes underperforming even compared to the no-retriever baseline, such as on STaRK-Prime using Biomixtral. The "w/o retriever" baseline, representing an LLM without retrieval augmentation, performs surprisingly well in some settings, indi-

	Approach	llama 3.1 8b						llama2-chat-7b						Biomixtral 7b					
		Exact Match			Rouge-1			Exact Match			Rouge-1			Exact Match			Rouge-1		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
RiTeK-PharmKG	G-retriever (He et al., 2024)	<b>36.97</b>	<b>46.07</b>	<b>38.31</b>	33.04	47.03	38.41	38.71	37.11	37.62	39.78	39.18	39.31	<b>43.01</b>	<b>41.59</b>	<b>42.01</b>	<b>43.95</b>	<b>42.69</b>	<b>43.10</b>
	GNN-RAG	33.21	43.01	37.31	21.00	44.89	26.00	<b>50.78</b>	<b>49.28</b>	<b>49.72</b>	<b>51.66</b>	50.29	<b>50.73</b>	39.93	39.12	39.26	41.69	40.89	41.08
	w/o retriever	32.45	43.40	34.23	<b>47.60</b>	<b>46.59</b>	<b>46.84</b>	38.91	37.63	38.02	40.57	39.31	39.72	41.49	39.43	40.05	41.25	41.06	40.99
RiTeK-ADint	G-retriever (He et al., 2024)	<b>50.83</b>	<b>50.07</b>	<b>50.31</b>	<b>57.34</b>	<b>56.61</b>	<b>56.87</b>	47.93	47.16	47.41	54.68	54.00	54.24	48.34	47.48	47.75	55.58	54.73	55.02
	GNN-RAG	40.88	40.90	40.43	44.43	45.01	45.49	<b>51.04</b>	<b>50.59</b>	<b>50.55</b>	<b>56.49</b>	<b>56.09</b>	<b>56.09</b>	<b>50.83</b>	<b>50.07</b>	<b>50.31</b>	<b>57.34</b>	<b>56.61</b>	<b>56.87</b>
	w/o retriever	49.59	48.48	48.82	55.23	54.29	54.61	46.58	45.82	46.06	51.66	49.91	46.47	49.79	48.93	49.20	56.43	53.80	54.15
STaRK-Prime	G-retriever (He et al., 2024)	<b>16.14</b>	<b>16.47</b>	<b>14.11</b>	17.21	27.86	19.21	10.15	8.45	8.17	21.75	18.08	18.40	12.22	<b>11.54</b>	10.58	23.15	21.37	20.72
	GNN-RAG	7.81	16.67	9.35	<b>18.13</b>	<b>27.50</b>	<b>19.65</b>	<b>16.00</b>	<b>15.04</b>	<b>14.50</b>	<b>24.78</b>	<b>23.51</b>	<b>22.99</b>	11.20	10.31	10.65	17.98	18.09	18.32
	w/o retriever	12.96	14.99	11.91	16.80	25.75	18.41	11.77	10.38	9.68	20.65	21.59	17.66	<b>12.96</b>	11.12	<b>10.73</b>	<b>24.83</b>	<b>21.59</b>	<b>21.68</b>

Table 5: Performance of different retrieval models across backbone LLMs. “w/o retriever” denotes an LLM without retrieval augmentation.

<p><b>Question:</b> How does Small RNA affect Organism Function that leads to Pathologic Function involving damaged myocardium?</p> <p><b>Ground Truth Path:</b> Small RNA → affects → Immune response → affects → Myocardial Reperfusion Injury</p> <p><b>Answers:</b> Myocardial Reperfusion Injury</p>	
G-retriever	<p><b>Path:</b> Small RNA → CAUSES → Fibrosis, MicroRNAs → ASSOCIATED_WITH → Myocardial Reperfusion Injury, Functional RNA → AFFECTS → biological adaptation to stress</p> <p><b>Answer:</b> Myocardial Reperfusion Injury</p>
GCR	<p><b>Answer:</b> Myocardial Reperfusion Injury, Myocardial Infarction</p>
GNN-RAG	<p><b>Path:</b> Small RNA → AFFECTS → Immune response, Immune response → AFFECTS → Myocardial Reperfusion Injury, Small RNA → AFFECTS → Cell physiology, Small RNA → AFFECTS → Virulence</p> <p><b>Answer:</b> Myocardial Reperfusion Injury</p>
<p><b>Question:</b> What disease is associated with the ch31l mutations affecting biomarkers linked to npy, characterized by overexpression in alb and known for leading to regressive behavior?</p> <p><b>Ground Truth Path:</b> ch31l → Mutations affect, or polymorphisms alter risk → schizophrenia → Biomarkers (diagnostic), or regulation linked to disease → npy, schizophrenia → Overexpression in disease → alb</p> <p><b>Answers:</b> Schizophrenia</p>	
TOG	<p><b>Path:</b> CH31L → AFFECTS → tube morphogenesis, CH31L → PREDISPOSES → Alzheimer’s Disease, CH31L → ASSOCIATED_WITH → Alzheimer’s Disease, Alzheimer’s Disease → AFFECTS → Abnormal behavior, Alzheimer’s Disease → ASSOCIATED_WITH → Gene Mutation, Alzheimer’s Disease → AFFECTS → Neurodegenerative Disorders, Neurodegenerative Disorders → COEXISTES_WITH → Abnormal degeneration, Neurodegenerative Disorders → COEXISTES_WITH → Death (finding), Neurodegenerative Disorders → COEXISTES_WITH → Impaired cognition</p> <p><b>Answer:</b> Alzheimer’s Disease, Neurodegenerative Disorders</p>
MCTS	<p><b>Path:</b> CH31L → PREDISPOSES → Alzheimer’s Disease → COEXISTES_WITH → Elevated serum protein → COEXISTES_WITH → Juvenile Huntington Disease → COEXISTES_WITH → Elevated serum protein</p> <p><b>Answer:</b> Juvenile Huntington Disease</p>
GPT-4	<p><b>Answer:</b> Autism Spectrum Disorder (ASD)</p>

Figure 2: A case study on RiTeK

579 cating that strong LLMs alone can capture a significant amount of relevant knowledge. For example, on RiTeK-ADint with Biomixtral, it achieves a Rouge-1 F1 score of 54.15, close to the GNN-RAG. However, in most cases, retrieval-augmented methods still yield superior performance. Notably, Biomixtral 7b tends to outperform the other two LLMs when combined with retrieval, especially in recall and F1. These results suggest that both the choice of retrieval strategy and the backbone LLM significantly impact end-task performance.

### 5.3.2 Case Study of Path and Answer Quality

591 We conduct a qualitative analysis to compare the reasoning paths and predicted answers from different retrieval models on two biomedical question-answering examples. As shown in Figure 2, all models successfully predicted the correct answer *Myocardial Reperfusion Injury* in the first case, although their reasoning paths varied in granularity and relevance. G-RETRIEVER and GNN-RAG

599 produced informative multi-hop paths that partially overlapped with the ground truth. 600

601 In contrast, for the second question involving *CH31L* and schizophrenia, only the ground truth path led to the correct answer. All baseline models failed: TOG and MCTS generated incorrect reasoning chains centered around *Alzheimer’s Disease* and *Juvenile Huntington Disease*, while GPT-4 hallucinated *Autism Spectrum Disorder*. These errors reveal the challenge of modeling rare or indirect biomedical associations, especially when entity relations involve subtle phenotypic markers. This case highlights the importance of precise multi-hop reasoning and clinically aligned retrieval in semi-structured biomedical graphs. 610 611 612 613

## 6 Conclusion

614 We present RiTeK, the first dataset specifically designed to evaluate the capability of models in handling complex reasoning over textual knowledge graphs (TKGs). This dataset offers diverse topological structures, relational types, entity types, and queries that integrate relational and textual information, requiring sophisticated reasoning across TKGs. RiTeK also includes rich textual descriptions for each node. To ensure the authenticity and accuracy of the queries, medical experts performed stringent validation. RiTeK sets a new standard for evaluating real-world retrieval systems. We evaluated 11 retrieval models on our benchmark dataset. Our experiments on RiTeK reveal significant challenges faced by current models in effectively handling both textual and relational information, especially under complex topological structures involving intricate relations and entities. RiTeK paves the way for future research aimed at advancing retrieval systems by emphasizing the need to enhance reasoning capabilities, particularly in retrieving complex reasoning paths under answer attribute constraints. 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637

## 7 Limitations

RiTeK is currently limited to queries that involve only a single topic entity and rely solely on the textual and structural information present in the graph. Future work should explore the inclusion of multiple topic entities and incorporate additional modalities, such as images, to enable a more comprehensive and robust information retrieval system.

Although we employed four medical experts for human evaluation, increasing the number of qualified domain experts would improve the statistical significance and robustness of our findings. Future work should consider expanding the pool of experts and addressing issues of fairness, and potential biases inherent in LLMs.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Junwei Bao, Nan Duan, Zhao Yan, Ming Zhou, and Tiejun Zhao. 2016. Constraint-based question answering with knowledge graph. In *Proceedings of COLING 2016, the 26th international conference on computational linguistics: technical papers*, pages 2503–2514.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.
- Guillaume Maurice Jean-Bernard Chaslot Chaslot. 2010. Monte-carlo tree search.
- Arman Cohan and Nazli Goharian. 2016. Revisiting summarization evaluation for scientific articles. *arXiv preprint arXiv:1604.00400*.
- Matthew Dunn, Levent Sagun, Mike Higgins, V Ugur Guney, Volkan Cirik, and Kyunghyun Cho. 2017. Searchqa: A new q&a dataset augmented with context from a search engine. *arXiv preprint arXiv:1704.05179*.
- Yu Gu, Sue Kase, Michelle Vanni, Brian Sadler, Percy Liang, Xifeng Yan, and Yu Su. 2021. Beyond iid: three levels of generalization for question answering on knowledge bases. In *Proceedings of the Web Conference 2021*, pages 3477–3488.
- Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and

- Bryan Hooi. 2024. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *arXiv preprint arXiv:2402.07630*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.
- Mingchen Li and Lifu Huang. 2023. Understand the dynamic world: An end-to-end knowledge informed framework for open domain entity state tracking. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 842–851.
- Mingchen Li and Shihao Ji. 2022. Semantic structure based query graph prediction for question answering over knowledge graph. *arXiv preprint arXiv:2204.10194*.
- Xinna Lin, Siqi Ma, Junjie Shan, Xiaojing Zhang, Shell Xu Hu, Tiannan Guo, Stan Z. Li, and Kaicheng Yu. 2024. [Biokgbench: A knowledge graph checking benchmark of ai agent for biomedical science](#). *Preprint*, arXiv:2407.00466.
- László Lovász. 1993. Random walks on graphs: A survey. In *Combinatorics, Paul Erdős is Eighty*, volume 2, pages 1–46. János Bolyai Mathematical Society.
- Linhao Luo, Zicheng Zhao, Chen Gong, Gholamreza Haffari, and Shirui Pan. 2024. Graph-constrained reasoning: Faithful reasoning on knowledge graphs with large language models. *arXiv preprint arXiv:2410.13080*.
- Costas Mavromatis and George Karypis. 2024. Gnnrag: Graph neural retrieval for large language model reasoning. *arXiv preprint arXiv:2405.20139*.
- P Rajpurkar. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Jiashuo Sun, Chengjin Xu, Luminyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung Shum, and Jian Guo. 2023. Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph. *arXiv preprint arXiv:2307.07697*.

744	Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris,	2024a. Tree of thoughts: Deliberate problem solving	800
745	Alessandro Sordoni, Philip Bachman, and Kaheer	with large language models. <i>Advances in Neural</i>	801
746	Suleman. 2016. Newsqa: A machine comprehension	<i>Information Processing Systems</i> , 36.	802
747			
748	Priyansh Trivedi, Gaurav Maheshwari, Mohnish Dubey,	Zonghai Yao, Zihao Zhang, Chaolong Tang, Xingyu	803
749	and Jens Lehmann. 2017. Lc-quad: A corpus for	Bian, Youxia Zhao, Zhichao Yang, Junda Wang,	804
750	complex question answering over knowledge graphs.	Huixue Zhou, Won Seok Jang, Feiyun Ouyang, et al.	805
751	In <i>The Semantic Web–ISWC 2017: 16th International</i>	2024b. Medqa-cs: Benchmarking large language	806
752	<i>Semantic Web Conference, Vienna, Austria, October</i>	models clinical skills using an ai-sce framework.	807
753	<i>21-25, 2017, Proceedings, Part II 16</i> , pages 210–218.	<i>arXiv preprint arXiv:2410.01553</i> .	808
754	Springer.		
755	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	Wen Tau Yih, Matthew Richardson, Chris Meek,	809
756	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	Ming Wei Chang, and Jina Suh. 2016. The value	810
757	et al. 2022. Chain-of-thought prompting elicits rea-	of semantic parse labeling for knowledge base ques-	811
758	soning in large language models. <i>Advances in neural</i>	tion answering. In <i>Proceedings of the 54th Annual</i>	812
759	<i>information processing systems</i> , 35:24824–24837.	<i>Meeting of the Association for Computational Lin-</i>	813
		<i>guistics (Volume 2: Short Papers)</i> .	814
760	Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel.	Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexan-	815
761	2018. Constructing datasets for multi-hop reading	der J Smola, and Le Song. 2018. Variational reason-	816
762	comprehension across documents. <i>Transactions of</i>	ing for question answering with knowledge graph. In	817
763	<i>the Association for Computational Linguistics</i> , 6:287–	<i>Thirty-Second AAAI Conference on Artificial Intelli-</i>	818
764	302.	<i>gence</i> .	819
765	Shirley Wu, Shiyu Zhao, Qian Huang, Kexin Huang,	Shuangjia Zheng, Jiahua Rao, Ying Song, Jixian Zhang,	820
766	Michihiro Yasunaga, Kaidi Cao, Vassilis N Ioan-	Xianglu Xiao, Evandro Fei Fang, Yuedong Yang, and	821
767	nidis, Karthik Subbian, Jure Leskovec, and James	Zhangming Niu. 2021. Pharmkg: a dedicated knowl-	822
768	Zou. 2024a. Avatar: Optimizing llm agents for	edge graph benchmark for biomedical data mining.	823
769	tool-assisted knowledge retrieval. <i>arXiv preprint</i>	<i>Briefings in bioinformatics</i> , 22(4):bbaa344.	824
770	<i>arXiv:2406.11200</i> .		
771	Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin	<b>A Appendix</b>	825
772	Huang, Kaidi Cao, Qian Huang, Vassilis N Ioanni-	<b>A.1 Ethics Statement</b>	826
773	dis, Karthik Subbian, James Zou, and Jure Leskovec.	All experiments in this study were conducted us-	827
774	2024b. Stark: Benchmarking llm retrieval on text-	ing publicly available datasets, including Prime,	828
775	tual and relational knowledge bases. <i>arXiv preprint</i>	ADint, and PharmKG. These datasets contain only	829
776	<i>arXiv:2404.13207</i> .	de-identified, non-personal, and non-sensitive in-	830
777	Yu Xia, Junda Wu, Sungchul Kim, Tong Yu, Ryan A	formation that was released for research purposes	831
778	Rossi, Haoliang Wang, and Julian McAuley. 2024.	under appropriate data licenses. No private or con-	832
779	Knowledge-aware query expansion with large lan-	fidential patient data were accessed or used. Con-	833
780	guage models for textual and relational retrieval.	sequently, no additional ethical approval was re-	834
781	<i>arXiv preprint arXiv:2410.13765</i> .	quired. We employed AI tools to assist with gram-	835
782	Yongkang Xiao, Yu Hou, Huixue Zhou, Gayo Di-	mar revision.	836
783	alloy, Marcelo Fiszman, Julian Wolfson, Li Zhou,	<b>A.2 Hyperparameters</b>	837
784	Halil Kilicoglu, You Chen, Chang Su, et al. 2024.	In the retrieval model, we set the maximum search	838
785	Repurposing non-pharmacological interventions for	depth to $d = 3$ for both Random Walk and MCTS.	839
786	alzheimer’s disease through link prediction on	The number of rollouts in MCTS is 16. For G-	840
787	biomedical literature. <i>Scientific reports</i> , 14(1):8693.	Retriever, KAR, GNN-RAG, and GCR, we adopted	841
788	Xi Yan, Patrick Westphal, Jan Seliger, and Ricardo Us-	the same hyperparameters as reported in their re-	842
789	beck. 2024. Bridging the gap: Generating a compre-	spective source papers.	843
790	hensive biomedical knowledge graph question an-	<b>A.3 TKG resources</b>	844
791	swering dataset. In <i>ECAI 2024</i> , pages 1198–1205.	Ensembl <sup>4</sup> , UMLS <sup>5</sup> , and Mondo Disease Ontol-	845
792	IOS Press.	ogy <sup>6</sup> .	846
793	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-		
794	gio, William W Cohen, Ruslan Salakhutdinov, and		
795	Christopher D Manning. 2018. Hotpotqa: A dataset		
796	for diverse, explainable multi-hop question answer-		
797	ing. <i>arXiv preprint arXiv:1809.09600</i> .		
798	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,		
799	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.		

<sup>4</sup><https://useast.ensembl.org/index.html>

<sup>5</sup><https://www.nlm.nih.gov/research/umls/index.html>

<sup>6</sup><https://mondo.monarchinitiative.org/>

847	<b>A.4 Medical textual knowledge graph</b>	textual properties and relational templates. Write	897
848	<b>construction</b>	the query from the perspective of a <persona>, en-	898
849	We construct two medical TKGs based on	sureing it is concise, human-like, and paraphrased	899
850	PharmKG (Zheng et al., 2021) and ADInt (Xiao	while retaining the original meaning.	900
851	et al., 2024), as the increased number of entity and	Consider the following characteristics for the	901
852	relation types introduces significant challenges for	persona:	902
853	path retrieval in the question answering over textual		
854	knowledge graph. We present the statistics of the	• Doctor: Formulate direct and practical ques-	903
855	relational structure in Table 1 and introduce each	tions aimed at diagnosing and treating. These	904
856	TKG as follows:	questions should focus on side effects, symp-	905
857	<b>PharmKG Textual Knowledge Graph:</b> We	toms, complications, and other clinically rele-	906
858	leverage the existing medical knowledge graph	vant aspects.	907
859	PharmKG (Zheng et al., 2021) which is a multi-	• Medical Scientist: Generate detailed and spe-	908
860	relational, attribute-rich biomedical knowledge	cific questions reflecting the complexity of	909
861	graph (KG) constructed from six publicly available	scientific inquiry. These questions should ex-	910
862	databases that provide high-quality structured in-	plorer etiology, pathophysiology, genetic fac-	911
863	formation. These databases include OMIM, Drug-	tors, pathways, proteins, or molecular func-	912
864	Bank, PharmGKB, Therapeutic Target Database	tions.	913
865	(TTD), SIDER, and HumanNet. PharmKG consists		
866	of over 500,000 distinct interconnections between	• Patient: Create straightforward questions	914
867	genes, drugs, and diseases, encompassing 29 types	that avoid professional medical terminology.	915
868	of relationships within a vocabulary of approxi-	These questions should focus on practical con-	916
869	mately 8,000 disambiguated entities. To enhance	cerns, such as symptoms, effects, inheritance,	917
870	the entity attributes, we incorporate textual details	or other relatable aspects, and may include	918
871	from various databases, including Ensembl, UMLS,	more context from daily life.	919
872	and Mondo Disease Ontology, as supplementary		
873	data sources.	<b>Textual Properties:</b> [ <input< td=""> <td>920</td> </input<>	920
874	<b>ADInt Textual Knowledge Graph:</b>	textual_properties>]	921
875	ADInt(Xiao et al., 2024) is a comprehensive	<b>Relational Templates:</b> [ <input< td=""> <td>922</td> </input<>	922
876	knowledge graph (KG) constructed from biomed-	relational_templates>]	923
877	ical literature, focusing on non-pharmacological	<b>Persona:</b> <input_persona> (e.g., Doctor, Medi-	924
878	interventions (NPI) and their associations with	cal Scientist, Patient)	925
879	Alzheimer’s disease (AD). ADInt includes 162,212	Ensure the query is realistic and diverse, lever-	926
880	entities spanning 113 UMLS semantic types,	aging flexibility in how the textual and relational	927
881	which, upon further classification, consist of	elements are presented. Avoid directly copying the	928
882	25,604 drugs, 16,474 diseases, 46,060 genes	input phrases; instead, paraphrase them while re-	929
883	and proteins, 2,525 dietary supplements (DS),	taining their original meaning. Please output only	930
884	and 128 complementary and integrative health	the generated query without any additional com-	931
885	(CIH) interventions. Moreover, ADInt contains	ments or explanations.	932
886	1,017,284 triples, capturing 15 distinct relation		
887	types, offering a rich dataset for exploring the intri-	<b>A.6 Data Analysis of query length and answer</b>	933
888	cate relationships between NPIs and AD. Same as	<b>length</b>	934
889	PharmKG, we also incorporate textual details from	We analyzed the distribution of query lengths (i.e.,	935
890	various databases, including Ensembl, UMLS, and	the number of words in each query) to assess the	936
891	Mondo Disease Ontology, as supplementary data	complexity of the queries and the amount of in-	937
892	sources.	formation they contain. As shown in the Figure 3,	938
893	<b>A.5 The prompt of Combining Textual and</b>	the query lengths range from 5 to 40 words, with	939
894	<b>Relational Information</b>	approximately 69% and 61% of queries in the two	940
895	You are a creative assistant tasked with generating	datasets having lengths between 15 and 25 words.	941
896	natural, diverse, and realistic queries by combining	Then, we analyzed the proportion of ground truth	942
		answers associated with each query. Generally,	943
		the more ground truth answers there are, the less	944

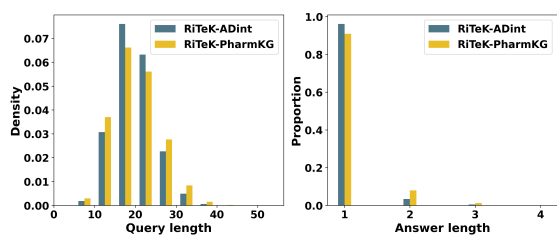


Figure 3: Distribution of query lengths and answer lengths on RiTeK-ADint and RiTeK-PharmKG datasets

precise the textual requirements in the query tend to be. To increase the difficulty of the question-answering task, we filtered out queries with too many ground truth answers during the dataset creation process, retaining only those with a maximum of three ground truth answers. In both datasets, over 90% of queries have a single ground truth answer, indicating that our queries are enriched with detailed textual information from entity attributes. This introduces more challenges when developing new graph retrieval methods

## B Relational Template

### B.1 RiTeK-PharmKG

1. Gene -> [Production by cell population] -> Gene
2. Gene -> [Enhance response, or activate, stimulate] -> Gene
3. Gene -> [Relationships involving regulation and pathways] -> Gene
4. Gene -> [Binding, ligand] -> Gene
5. Gene -> [Affects expression/production] -> Gene
6. Gene -> [Gene-Gene] -> Gene
7. Chemical -> [Chemical-Chemical] -> Chemical
8. Disease -> [Ancestors of disease] -> Disease
9. Disease -> [Associations between diseases] -> Disease
10. Gene -> [Interactions] -> Chemical
11. Chemical -> [Interactions] -> Gene
12. Gene -> [Interactions] -> Gene
13. Gene -> [Interactions] -> Disease
14. Gene -> [Drug targets] -> Disease

15. Gene -> [Role in pathogenesis, or promotes progression] -> Disease 978  
979
16. Gene -> [Mutations affect, or polymorphisms alter risk] -> Disease 980  
981
17. Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene 982  
983
18. Disease -> [Overexpression in disease] -> Gene 984  
985
19. Chemical -> [Treatment or therapy] -> Disease 986  
987
20. Chemical -> [Side effect or adverse event] -> Disease 988  
989
21. Chemical -> [Inhibits cell growth] -> Disease 990
22. Chemical -> [Role in pathogenesis] -> Disease 991  
992
23. Chemical -> [Prevents, suppresses, or alleviates, reduces] -> Disease 993  
994
24. Disease -> [Biomarkers (progression)] -> Chemical 995  
996
25. Chemical -> [Agonism, activation, or antagonism, blocking] -> Gene 997  
998
26. Chemical -> [Binding, ligand] -> Gene 999
27. Chemical -> [Affects expression/production] -> Gene 1000  
1001
28. Chemical -> [Inhibits] -> Gene 1002
29. Gene -> [Transport, channels] -> Chemical 1003
30. Gene -> [Metabolism, pharmacokinetics] -> Chemical 1004  
1005
31. Gene -> [Enzyme activity] -> Chemical 1006
32. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Drug targets] -> Disease 1007  
1008
33. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Role in pathogenesis, or promotes progression] -> Disease 1009  
1010
34. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Mutations affect, or polymorphisms alter risk] -> Disease 1012  
1013
35. Gene -> [Relationships involving regulation and pathways] -> Gene -> [Binding, ligand] -> Gene 1015  
1016  
1017

1018	36. Gene -> [Binding, ligand] -> Gene -> [Affects expression/production] -> Gene	52. Gene -> [Metabolism, pharmacokinetics] -> Chemical -> [Binding, ligand] -> Gene -> [Inhibits] -> Chemical	1061
1019			1062
1020	37. Gene -> [Interactions] -> Gene -> [Interactions] -> Chemical	53. Gene -> [Interactions] -> Chemical -> [Treatment or therapy] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene	1063
1021			1064
1022	38. Gene -> [Transport, channels] -> Chemical -> [Agonism, activation, or antagonism, blocking] -> Gene		1065
1023			1066
1024			1067
1025	39. Gene -> [Metabolism, pharmacokinetics] -> Chemical -> [Binding, ligand] -> Gene	54. Gene -> [Interactions] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene -> [Transport, channels] -> Chemical	1068
1026			1069
1027	40. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Enhance response, or activate, stimulate] -> Gene		1070
1028			1071
1029		55. Gene -> [Role in pathogenesis, or promotes progression] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene -> [Metabolism, pharmacokinetics] -> Chemical	1072
1030	41. Gene -> [Interactions] -> Chemical -> [Treatment or therapy] -> Disease		1073
1031			1074
1032	42. Gene -> [Interactions] -> Chemical -> [Side effect or adverse event] -> Disease		1075
1033			1076
1034	43. Gene -> [Interactions] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene	56. Chemical -> [Agonism, activation, or antagonism, blocking] -> Gene -> [Drug targets] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene	1077
1035			1078
1036			1079
1037	44. Chemical -> [Treatment or therapy] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene	57. Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene -> [Role in pathogenesis, or promotes progression] -> Disease	1080
1038			1081
1039			1082
1040	45. Disease -> [Associations between diseases] -> Disease -> [Ancestors of disease] -> Disease		1083
1041			1084
1042	46. Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene	58. Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene -> [Metabolism, pharmacokinetics] -> Chemical -> [Side effect or adverse event] -> Disease	1085
1043			1086
1044			1087
1045			1088
1046	47. Gene -> [Interactions] -> Gene -> [Transport, channels] -> Chemical	59. Gene -> [Production by cell population] -> Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Relationships involving regulation and pathways] -> Gene	1089
1047			1090
1048	48. Gene -> [Metabolism, pharmacokinetics] -> Chemical -> [Binding, ligand] -> Gene		1091
1049			1092
1050	49. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Drug targets] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene	60. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Binding, ligand] -> Gene -> [Affects expression/production] -> Gene	1093
1051			1094
1052			1095
1053			1096
1054	50. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Mutations affect, or polymorphisms alter risk] -> Disease -> [Overexpression in disease] -> Gene	61. Gene -> [Relationships involving regulation and pathways] -> Gene -> [Gene-Gene] -> Gene -> [Binding, ligand] -> Gene	1097
1055			1098
1056			1099
1057			1100
1058	51. Gene -> [Transport, channels] -> Chemical -> [Agonism, activation, or antagonism, blocking] -> Gene -> [Binding, ligand] -> Chemical	62. Gene -> [Interactions] -> Gene -> [Interactions] -> Gene -> [Transport, channels] -> Chemical	1101
1059			1102
1060		63. Gene -> [Interactions] -> Gene -> [Interactions] -> Gene -> [Metabolism, pharmacokinetics] -> Chemical	1103
			1104
			1105

1106	64. Gene -> [Enhance response, or activate, stimulate] -> Gene -> [Mutations affect, or polymorphisms alter risk] -> Disease -> [Overexpression in disease] -> Gene	14. Biologically Active Substance -> [affects] -> Organism Function	1147
1107			1148
1108		15. Biologically Active Substance -> [causes] -> Injury or Poisoning	1149
1109			1150
1110	65. Gene -> [Enzyme activity] -> Chemical -> [Affects expression/production] -> Gene -> [Chemical-Chemical] -> Chemical	16. Biologically Active Substance -> [disrupts] -> Gene or Genome	1151
1111			1152
1112			
1113	66. Gene -> [Interactions] -> Chemical -> [Role in pathogenesis] -> Disease -> [Overexpression in disease] -> Gene	17. Body Part, Organ, or Organ Component -> [produces] -> Immunologic Factor	1153
1114			1154
1115		18. Cell Component -> [affects] -> Molecular Function	1155
1116	67. Chemical -> [Side effect or adverse event] -> Disease -> [Biomarkers (diagnostic), or regulation linked to disease] -> Gene -> [Mutations affect, or polymorphisms alter risk] -> Disease		1156
1117		19. Cell Component -> [produces] -> Nucleic Acid, Nucleoside, or Nucleotide	1157
1118			1158
1119		20. Cell Function -> [affects] -> Mental or Behavioral Dysfunction	1159
1120			1160
1121	68. Chemical -> [Inhibits cell growth] -> Disease -> [Overexpression in disease] -> Gene -> [Role in pathogenesis, or promotes progression] -> Disease	21. Cell Function -> [produces] -> Biologically Active Substance	1161
1122			1162
1123		22. Cell or Molecular Dysfunction -> [affects] -> Neoplastic Process	1163
1124			1164
1125	<b>B.2 RiTeK-ADint</b>		
1126	1. Amino Acid, Peptide, or Protein -> [affects] -> Cell Function	23. Cell or Molecular Dysfunction -> [manifestation of] -> Pathologic Function	1165
1127			1166
1128	2. Amino Acid, Peptide, or Protein -> [affects] -> Disease or Syndrome	24. Cell -> [produces] -> Organic Chemical	1167
1129			
1130	3. Amino Acid, Peptide, or Protein -> [causes] -> Anatomical Abnormality	25. Congenital Abnormality -> [affects] -> Virus	1168
1131		26. Congenital Abnormality -> [manifestation of] -> Organism Function	1169
1132	4. Amino Acid, Peptide, or Protein -> [interacts with] -> Pharmacologic Substance		1170
1133		27. Diagnostic Procedure -> [affects] -> Genetic Function	1171
1134	5. Anatomical Abnormality -> [affects] -> Organ or Tissue Function		1172
1135		28. Disease or Syndrome -> [affects] -> Organ or Tissue Function	1173
1136	6. Anatomical Abnormality -> [complicates] -> Disease or Syndrome		1174
1137		29. Disease or Syndrome -> [associated with] -> Therapeutic or Preventive Procedure	1175
1138	7. Anatomical Abnormality -> [manifestation of] -> Genetic Function		1176
1139		30. Disease or Syndrome -> [manifestation of] -> Cell or Molecular Dysfunction	1177
1140	8. Antibiotic -> [affects] -> Molecular Function		1178
1141	9. Antibiotic -> [causes] -> Pathologic Function	31. Finding -> [manifestation of] -> Pathologic Function	1179
1142	10. Antibiotic -> [disrupts] -> Cell Component		1180
1143	11. Antibiotic -> [treats] -> Disease or Syndrome	32. Gene or Genome -> [produces] -> Amino Acid, Peptide, or Protein	1181
1144	12. Bacterium -> [causes] -> Cell or Molecular Dysfunction		1182
1145		33. Genetic Function -> [affects] -> Human	1183
1146	13. Bacterium -> [interacts with] -> Human	34. Genetic Function -> [produces] -> Cell Component	1184
			1185

1186	35. Hazardous or Poisonous Substance -> [affects] -> Mental or Behavioral Dysfunction	55. Pharmacologic Substance -> [treats] -> Sign or Symptom	1224
1187			1225
1188	36. Hazardous or Poisonous Substance -> [disrupts] -> Organ or Tissue Function	56. Sign or Symptom -> [manifestation of] -> Genetic Function	1226
1189			1227
1190	37. Health Care Activity -> [affects] -> Disease or Syndrome	57. Therapeutic or Preventive Procedure -> [affects] -> Neoplastic Process	1228
1191			1229
1192	38. Human -> [interacts with] -> Human	58. Virus -> [interacts with] -> Human	1230
1193	39. Immunologic Factor -> [affects] -> Pathologic Function	<b>C Topological Structures</b>	1231
1194			
1195	40. Indicator, Reagent, or Diagnostic Aid -> [interacts with] -> Hazardous or Poisonous Substance	<b>C.1 Definition of Six Reasoning Topologies</b>	1232
1196		Following the semantic structure framework introduced by (Li and Ji, 2022), we define six reasoning topologies that serve as the structural backbone for question generation and reasoning path simulation in RiTeK. Each topology represents a distinct reasoning pattern between entities in a textual knowledge graph (TKG), reflecting different levels of relational complexity and logical dependency.	1233
1197			1234
1198	41. Injury or Poisoning -> [disrupts] -> Genetic Function		1235
1199			1236
1200	42. Medical Device -> [treats] -> Mental or Behavioral Dysfunction		1237
1201			1238
1202	43. Mental or Behavioral Dysfunction -> [affects] -> Organism Function	• 1-hop: The simplest reasoning structure, consisting of a single relation connecting the topic entity and the answer entity.	1239
1203			1240
1204	44. Molecular Function -> [affects] -> Virus		1241
1205	45. Neoplastic Process -> [affects] -> Bacterium	• 2-hop: A two-step linear reasoning chain where the answer is connected to the topic entity through an intermediate entity.	1242
1206			1243
1207	46. Neoplastic Process -> [associated with] -> Neoplastic Process	• 3-hop: A longer reasoning chain with three relational steps, representing more complex dependency and multi-level inference.	1244
1208			1245
1209	47. Nucleic Acid, Nucleoside, or Nucleotide -> [interacts with] -> Immunologic Factor	• 1-hop with constraint: A single relational edge combined with a semantic or categorical constraint that filters valid answers.	1246
1210			1247
1211	48. Organ or Tissue Function -> [produces] -> Immunologic Factor	• 2-hop with constraint: A two-hop reasoning chain where the final answer is subject to an additional semantic or categorical restriction.	1248
1212			1249
1213	49. Organic Chemical -> [affects] -> Pathologic Function	• Two-to-one (Converging paths): Two distinct entities or relations converge on a common target entity. This topology reflects intersective reasoning, where the answer satisfies multiple relational constraints simultaneously.	1250
1214			1251
1215	50. Organic Chemical -> [interacts with] -> Pharmacologic Substance		1252
1216			1253
1217	51. Organism Function -> [affects] -> Disease or Syndrome		1254
1218			1255
1219	52. Pathologic Function -> [associated with] -> Therapeutic or Preventive Procedure		1256
1220			1257
1221	53. Pathologic Function -> [manifestation of] -> Organ or Tissue Function	<b>C.2 Examples in the Medical Domain</b>	1258
1222			1259
1223	54. Pharmacologic Substance -> [affects] -> Genetic Function	• 1-hop: For example, Amino Acid, Peptide, or Protein → [affects] → Disease or Syndrome illustrates a direct relationship such as “Which diseases are affected by a given protein?”	1260
			1261
			1262
			1263
			1264
			1265

	Shannon Entropy	Type-Token Ratio
Medical domain		
RiTeK-ADint	10.04	0.187
RiTeK-PharmKG	9.61	0.157
STARK-PRIME	9.63	0.143
General domain		
STARK-AMAZON	10.39	0.179
STARK-MAG	10.25	0.180

Table 6: Query diversity measurement

- 1266 • 2-hop: An example is Gene  $\rightarrow$  [Interactions]  
1267  $\rightarrow$  Chemical  $\rightarrow$  [Treatment or therapy]  $\rightarrow$  Dis-  
1268 ease, corresponding to the question “Which  
1269 diseases can be treated by chemicals that in-  
1270 teract with a specific gene?”
- 1271 • 3-hop: For instance, Gene  $\rightarrow$  [Interactions]  
1272  $\rightarrow$  Chemical  $\rightarrow$  [Treatment or therapy]  $\rightarrow$   
1273 Disease  $\rightarrow$  [Biomarkers (diagnostic), or regu-  
1274 lation linked to disease]  $\rightarrow$  Gene models com-  
1275 plex interactions such as “Which genes serve  
1276 as biomarkers for diseases treated by chemi-  
1277 cals interacting with a given gene?”
- 1278 • 1-hop with constraint: For example, Phar-  
1279 macologic Substance  $\rightarrow$  [treats]  $\rightarrow$  Disease  
1280 or Syndrome under the constraint “Disease  
1281 type = neurodegenerative” represents ques-  
1282 tions like “Which pharmacologic substances  
1283 treat neurodegenerative diseases?”
- 1284 • 2-hop with constraint: An example is Gene  $\rightarrow$   
1285 [Enhance response, or activate, stimulate]  $\rightarrow$   
1286 Gene  $\rightarrow$  [Drug targets]  $\rightarrow$  Disease with the  
1287 constraint “Disease subtype = inflammatory”,  
1288 corresponding to “Which genes activate other  
1289 genes that target inflammatory diseases?”
- 1290 • Two-to-one (Converging paths): For instance,  
1291 Gene 1  $\rightarrow$  [Interactions]  $\rightarrow$  Chemical  $\leftarrow$  [In-  
1292 teractions]  $\leftarrow$  Gene 2 captures the question  
1293 “Which chemicals interact with both Gene A  
1294 and Gene B?”.

## 1295 D Data Distribution Analysis