
Diagnosing and Addressing Pitfalls in KG-RAG Datasets: Toward More Reliable Benchmarking

Liangliang Zhang¹ Zhuorui Jiang² Hongliang Chi¹ Haoyang Chen¹ Mohammed Elkoumy¹
Fali Wang³ Qiong Wu⁴ Zhengyi Zhou⁴ Shirui Pan⁵ Suhang Wang³ Yao Ma¹

¹Rensselaer Polytechnic Institute ²University of Toronto ³Pennsylvania State University
⁴AT&T Chief Data Office ⁵Griffith University

{zhangl41, chih3, chenh29, elkoum, may13}@rpi.edu
zhuorui.jiang@mail.utoronto.ca, {qw6547, zz547k}@att.com
s.pan@griffith.edu.au, {fqw5095, szw494}@psu.edu

Abstract

1 Knowledge Graph Question Answering (KGQA) systems rely on high-quality
2 benchmarks to evaluate complex multi-hop reasoning. However, despite their
3 widespread use, popular datasets such as WebQSP and CWQ suffer from critical
4 quality issues, including inaccurate or incomplete ground-truth annotations,
5 poorly constructed questions that are ambiguous, trivial, or unanswerable, and
6 outdated or inconsistent knowledge. Through a manual audit of 16 popular KGQA
7 datasets—including WebQSP and CWQ—we find that the average factual correctness
8 rate is only 57%. To address these issues, we introduce KGQAGen, an LLM-in-
9 the-loop framework that systematically resolves these pitfalls. KGQAGen combines
10 structured knowledge grounding, LLM-guided generation, and symbolic verification
11 to produce challenging and verifiable QA instances. Using KGQAGen, we
12 construct KGQAGen-10k, a 10K-scale benchmark grounded in Wikidata, and evaluate
13 a diverse set of KG-RAG models. Experimental results demonstrate that
14 even state-of-the-art systems struggle on this benchmark, highlighting its ability
15 to expose limitations of existing models. Our findings advocate for more rigorous
16 benchmark construction and position KGQAGen as a scalable framework for
17 advancing KGQA evaluation¹.

18 1 Introduction

19 Knowledge graph-based retrieval-augmented generation (KG-RAG) systems combine symbolic
20 retrieval with generative reasoning, enabling question answering that requires both factual accuracy
21 and structured inference [47, 12, 28, 21, 76, 56]. As these systems gain attention in both academic
22 and industrial settings [15, 49, 77], benchmark datasets play a central role in measuring progress and
23 guiding model development [73, 58, 79, 20, 10]. However, despite their central role in evaluation,
24 little attention has been paid to the quality and reliability of these benchmarks themselves.

25 To better understand the limitations of existing benchmarks, we conduct a detailed manual inspection
26 of 16 publicly available KGQA datasets including the widely adopted WebQSP and CWQ [73, 58,
27 5, 54, 39, 79, 4, 51, 62, 14, 23, 25, 20, 48, 8, 10]. A detailed summary of these datasets, along
28 with our sampling and evaluation protocol, is provided in Appendix B. Our analysis reveals several
29 recurring issues that compromise their utility for evaluating KG-RAG systems. These include

¹KGQAGen: <https://github.com/liangliang6v6/KGQAGen>; KGQAGen-10k: <https://huggingface.co/datasets/lianglz/KGQAGen-10k>.

30 factually incorrect or outdated ground truth answers, and ambiguously phrased or trivially simple
31 questions. In particular, WebQSP and CWQ, which serve as the dominant² evaluation benchmarks in
32 recent KG-RAG research [72, 18, 71, 66, 78], exhibit serious deficiencies. Specifically, we found
33 only 52% of sampled WebQSP examples and 49.3% of CWQ examples to be factually correct. Across
34 all 16 datasets, the average correctness rate is just 57%, based on manual evaluation of over 1,000
35 question–answer pairs. A detailed summary of dataset quality is provided in Table 1. Additionally,
36 many datasets rely on rigid exact-match metrics that penalize semantically correct answers expressed
37 in alternative forms, further limiting their reliability [50, 45, 64].

38 Motivated by the issues, we propose KGQAGen, a framework for constructing high-quality benchmarks
39 for KG-RAG systems. KGQAGen grounds question generation in a large, up-to-date Wikidata [63] as
40 knowledge base and leverages a modular LLM-in-the-loop pipeline to ensure that each instance is
41 factually correct and linguistically well-formed. Key components of the framework include iterative
42 LLM-guided KG exploration and symbolic verification. Each question is grounded in a subgraph of
43 the knowledge graph, which is iteratively expanded from a seed entity to include richer relational
44 structures. This expansion enables the generation of more challenging and semantically complex
45 questions. To ensure both efficiency and relevance, an LLM guides the expansion process by selecting
46 informative entities and checking contextual sufficiency. Finally, symbolic verification using SPARQL
47 ensures that the generated answers are correct and fully supported by the knowledge base. Together,
48 these components enable the scalable generation of challenging, diverse, and reliable QA pairs.

49 We further use KGQAGen to generate a sample dataset KGQAGen-10k consisting of 10,787 QA pairs.
50 KGQAGen-10k serves as a case study for investigating the quality, diversity, and characteristics of
51 instances produced by the framework. Manual inspection of 300 samples reveals 96% factual accu-
52 racy. We analyze the dataset along several axes—including linguistic complexity, topic coverage
53 that KGQAGen produces well-scoped, multi-hop questions with diverse structure and clear answer
54 grounding. Moreover, we use KGQAGen-10k to benchmark a diverse set of models, including both
55 pure LLMs and KG-RAG approaches. Our evaluation shows the difficulty of KGQAGen-10k: even
56 SOTA models such as GPT-4 . 1 and recent KG-RAG systems such as GCR [34] and PoG [9] achieve
57 only moderate performance. These results validate that the benchmark can effectively expose limita-
58 tions in retrieval and reasoning, and demonstrate the utility of KGQAGen as a scalable, interpretable,
59 and diagnostic tool for developing more capable KG-RAG systems. While KGQAGen-10k serves
60 as a representative sample for our investigation, KGQAGen is fully modular and scalable, making it
61 well-suited for constructing large-scale benchmarks with minimal human oversight.

62 **Contributions.** To summarize, our work makes three key contributions: (1) a systematic manual
63 audit of 16 widely used KGQA datasets, uncovering critical quality and evaluation issues; (2) the
64 development of KGQAGen, a scalable LLM-guided framework for generating challenging, grounded,
65 and verifiable KGQA benchmarks; and (3) the construction of KGQAGen-10k, a 10K-scale dataset
66 used to analyze question characteristics and benchmark a diverse set of KG-RAG models, revealing
67 significant performance gaps and opportunities for future improvement.

68 2 Related Work

69 In this section, we briefly review prior work on KG-RAG and KGQA benchmark construction. We
70 focus on the most relevant developments; a detailed discussion is provided in Appendix A.

71 **KG-RAG Methods.** KG-RAG systems enhance large language models by incorporating structured
72 knowledge from KGs. Recent approaches such as RoG [33], GCR [34], and ToG [72] combine
73 symbolic retrieval with multi-hop reasoning capabilities. Other frameworks, including DeCAF [74],
74 DualR [29], and FRAG [18], focus on diverse strategies for integrating retrieval and generation.
75 These efforts aim to improve factual accuracy and interpretability with external KG.

76 **KGQA Benchmarks.** Traditional KGQA datasets like WebQSP mainly rely on Freebase [7], which
77 officially dumped since 2016. WebQuestions [6] and its successor WebQSP [73] generate dataset
78 manually by collecting questions through Google Suggest API [38] and having Amazon MTurk
79 workers annotate the ground truths. CWQ [58] extended this approach to handle more complex queries.
80 To improve coverage and linguistic variety, subsequent benchmarks shifted to more diverse knowl-

²Around 30 papers on KG-RAG released between 2022 and 2025 adopted WebQSP and/or CWQ for evaluation, highlighting their widespread use. A detailed breakdown is provided in Appendix B

Table 1: Systematic Audit of KGQA datasets.

Dataset	KG	Year	Sample / Total	Correctness (%)
WebQSP [73]	Freebase	2016	100 / 1639	52.00
CWQ [58]	Freebase	2018	300 / 3531	49.33
ComplexQuestions [5]	Freebase	2016	60 / 800	63.33
GraphQuestions [54]	Freebase	2016	60 / 2607	70.00
QALD-9 [39]	DBpedia	2018	60 / 150	61.67
MetaQA [79]	WikiMovies	2018	60 / 39093	20.00
SimpleDBpediaQA [4]	DBpedia	2018	60 / 8595	43.33
CSQA [51]	Wikidata	2018	60 / 27797	65.00
LC-QuAD 1.0&2.0 [62, 14]	DBpedia/Wikidata	2017/2019	60 / 7046	38.34
FreeBaseQA [23]	Freebase	2019	60 / 3996	98.67
CFQ [25]	Freebase	2020	60 / 239357	71.67
GrailQA [20]	Freebase	2020	60 / 13231	30.00
QALD-9-Plus [48]	DB/Wikidata	2022	60 / 136	63.33
KQAPro [8]	FB15k+Wikidata	2022	60 / 11797	66.67
Dynamic-KGQA [10]	YAGO	2025	60 / 40000	45.00

edge bases. LC-QuAD 1.0&2.0 [62, 14] used a hybrid approach combining SPARQL templates and human annotation over DBpedia and Wikidata. QALD-9 [39, 48] targeted multilingual QA evaluation, while CSQA [51] introduced conversational structures. KQAPro [8] emphasized compositional reasoning with program annotations. MetaQA [79] used a synthetic movie KG to assess multi-hop reasoning and robustness to paraphrasing. Other datasets addressed specific evaluation goals. GrailQA [20] focused on generalization—i.i.d., compositional, and zero-shot—using questions over Freebase. GeneralizableKGQA [24] extended this by re-splitting multiple datasets under shared evaluation settings. CBench [44] and SmartBench [42] analyzed datasets from the SPARQL structural complexity and linguistic diversity. Maestro [43] proposed a rule-based automatic construction framework, though its reliance on manually defined predicate rules limits its generality. More recently, scalable generation methods have emerged. CHATTY-Gen [40] introduced dialogue-style questions featuring coreference and ellipsis. Dynamic-KGQA [10] employed LLMs to adaptively generate QA instances from YAGO 4.5, but faced challenges from KG sparsity and hallucinated outputs. In summary, these existing efforts aim to enhance KGQA datasets from various perspectives—such as increasing linguistic diversity, enabling compositional reasoning, and improving scalability through automation. While valuable, few have systematically examined or addressed quality assurance. In contrast, our work focuses on factual correctness and verifiability.

3 Pitfalls of Existing KGQA Benchmarks

We identify two key issues with existing KGQA benchmarks for KG-RAG: (1) data quality problems, including inaccurate labels and trivial or ambiguous questions; and (2) rigid EM-based evaluation, which overlooks semantically correct but differently phrased answers.

3.1 Dataset Quality Issues

To assess the reliability of existing KGQA benchmarks, we manually inspected over 1,000 question-answer pairs sampled from 16 widely used datasets. A summary of findings is provided in Table 1. A description of the datasets and inspection protocol is included in appendix B. For each dataset, we randomly selected a representative subset of examples and evaluated them along three key dimensions: (1) factual correctness of the annotated answer, and (2) clarity and appropriateness of the question and (3) current evaluation of exact match shortcoming. We paid particular attention to the WebQSP and CWQ datasets, as they are the most dominant benchmarks in recent KGQA research. As shown in Table 1, for WebQSP and CWQ, we sampled 100 and 300 examples, respectively. For the remaining datasets, we uniformly sampled 60 examples. Overall, our inspection revealed substantial quality issues across many benchmarks. Notably, the most widely used WebQSP and CWQ datasets show serious quality issues. In WebQSP, only 52% of the sampled answers were judged correct, with quality issues including inaccurate answer and poor questions. Similarly, CWQ, which builds on WebQSP as its seed, achieved a correctness rate of just 49.33%, suffering from similar flaws along with additional issues due to increased question complexity. Beyond these two, we found that over half of the evaluated datasets had correctness rates below 60%, suggesting widespread challenges with annotation accuracy and question clarity. Next, we discuss the specific issues.

119 3.1.1 Inaccurate Ground Truth Answers.

120 A major issue across many KGQA datasets is the presence of inaccurate ground truth answers. Our
121 detailed inspection reveals that such errors arise from distinct sources, which we categorize below.

- 122 • **Incorrect annotations** refer to cases where the labeled answer fails to align with the question’s
123 intent. For instance, in WebQSP, the question “*Where did Andy Murray start playing tennis?*” is
124 incorrectly annotated with *2005*, a year rather than a location—demonstrating a mismatch between
125 the expected answer type and the provided label.
- 126 • **Outdated answers** occur when the annotated response reflects facts that were once accurate but are
127 no longer valid according to up-to-date knowledge. This is common for time-sensitive information,
128 such as political positions, affiliations, or current locations. For instance, the question from WebQSP
129 “*Who is the president of Peru now?*” lists *Ollanta Humala*, who was president from 2011-2016. This
130 leads to penalizing correct model predictions that reflect up-to-date knowledge.
- 131 • **Incomplete annotations** happen when a question has multiple valid answers, but only one or a
132 few are labeled as correct. This is common in open-ended or set-based questions. For instance, in
133 Dynamic-KGQA, the question “*Which American actress is known for her notable work in a film
134 where she also starred as an actor?*” is labeled with just *Lindsay Lohan*, even though others like
135 *Barbra Streisand* and *Angelina Jolie* clearly qualify. As a result, models that return equally correct
136 answers are unfairly penalized.

137 We provide additional examples of each type of issues in Appendix C.1.1. Although FreeBaseQA
138 stands out with a high answer correctness rate (98.7%), a closer inspection reveals that many of its
139 questions are overly simple and lack reasoning depth. We detail this issue in the next subsection.

140 3.1.2 Low-Quality or Ambiguous Questions

141 Another major limitation of existing KGQA datasets is the prevalence of low-quality or poorly
142 constructed questions. In our analysis, these issues were observed in nearly all datasets to varying
143 degrees. We categorize the common problems into three types.

- 144 • **Ambiguous phrasing** arises when questions lack sufficient context to uniquely identify a specific
145 entity or relation. For example, the WebQSP question “*What does George Wilson do for a living?*”
146 is inherently under-specified because it fails to distinguish which individual named George Wilson
147 is being referred to. Multiple well-known figures share this name — including a fictional character
148 from *The Great Gatsby*, a recurring comic strip character from *Dennis the Menace*, and several
149 professional athletes, making it hard to answer accurately without additional contextual clues.
- 150 • **Low-complexity questions** require only shallow, one-hop retrieval, or string matching, offering
151 little value in evaluating complex reasoning. This issue is especially common in MetaQA and
152 FreeBaseQA. While MetaQA suffers from both low answer correctness (25%) and poor clarity,
153 FreeBaseQA achieves a high correctness rate (98.7%). However, this high correctness appears to
154 stem from the simplicity of the questions. To investigate this further, we used GPT-4o to answer
155 directly the questions of FreeBaseQA, achieving 90.39% accuracy in the evaluation of exact
156 matches. The strong performance even without KG access suggests that the dataset contains mostly
157 factoid, low-reasoning questions that are easily handled by powerful language models. These results
158 reinforce the view that FreeBaseQA, despite its clean annotations, does not present a sufficient
159 reasoning challenge for benchmarking KG-RAG systems.
- 160 • **Unanswerable, subjective, or ill-formed questions** are incompatible with structured KG-based
161 reasoning. For example, WebQSP includes questions such as “*What to do today in Atlanta with
162 kids?*” and “*What inspired Michael Jackson to become a singer?*”, both of which are inherently
163 subjective and lack definitive, factual answers. Similarly, the Dynamic-KGQA question “*Which
164 American university alumnus shares the same nationality as the creator of Kryptos?*” is poorly
165 constructed, as any American university alumnus would automatically satisfy the condition.

166 We provide additional representative examples for each issue type in Appendix C.1.2.

167 3.2 Limitations of Exact-Match Evaluation

168 Existing KGQA benchmarks also suffer from shortcomings in their evaluation protocols. Most
169 benchmarks rely on rigid exact-match criteria that fail to account for semantically correct answers
170 expressed in different surface forms. This leads to false negatives when models generate correct
171 answers that differ slightly from the annotated label. For example, in the WebQSP question “*What*

172 *is the Australian dollar called?*”, the ground truth answer is “AUD”. A prediction of “Australian
173 dollar”, though semantically equivalent, would be considered incorrect under exact match. Similar
174 mismatches arise from differences in formatting, paraphrasing, or entity naming conventions (e.g.,
175 “Germany” vs. “Federal Republic of Germany”). Additional examples are provided in Appendix C.2.

176 4 KGQAGen: A Framework for Grounded KGQA Dataset Construction

177 As discussed in the previous section, existing KGQA benchmarks suffer from quality and reliability
178 issues that limit their utility for evaluating KG-RAG systems. To address this, we introduce KGQAGen,
179 a modular framework for constructing high-quality QA datasets that are both semantically grounded
180 and verifiable. KGQAGen grounds each question in explicit KG evidence and leverages LLMs to assist
181 with subgraph expansion, question generation, and answer validation—enabling the scalable creation
182 of challenging and reliable benchmark instances.

183 The overall process of KGQAGen is illustrated in Figure 1. The framework consists of three main stages:
184 (1) **Seed Subgraph Initialization:** The process begins by selecting a seed entity and constructing a
185 local subgraph by retrieving related facts from the KG. This subgraph provides the initial context for
186 reasoning. (2) **Question Generation through Iterative LLM-Guided Subgraph Expansion:** To
187 support non-trivial, multi-hop questions, we iteratively expand the subgraph by traversing neighboring
188 entities and relations. This involves alternating between KG traversal and LLM evaluation. At each
189 step, the subgraph is expanded by exploring neighboring entities and relations within the KG. After
190 each expansion, an LLM is prompted to evaluate whether the subgraph contains sufficient information
191 to support a well-formed, multi-hop question. If not, the expansion continues. Once the subgraph
192 is judged sufficient, the LLM generates a natural language question, identifies the corresponding
193 answer set, extracts a minimal *supporting subgraph*, and constructs the associated SPARQL query. (3)
194 **Answer Validation and Refinement:** The final stage ensures that the generated answer set is correct
195 and fully grounded in the KG. To achieve this, the generated SPARQL query is executed against
196 the knowledge base to retrieve the actual answers. If the results match the generated answer set, the
197 instance is accepted. An ablation study presented in Appendix E further confirms the importance of
198 each component in KGQAGen. Both LLM-guided subgraph expansion and SPARQL-based verification
199 are shown to be essential for producing accurate and well-grounded QA pairs.

200 Before we proceed to detail the framework, we introduce the notation used throughout this section.
201 **Notations.** We denote a KG as a set of triples $\mathcal{G} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$, where \mathcal{E} is the set of entities and
202 \mathcal{R} is the set of relations. Each triple $\langle s, p, o \rangle \in \mathcal{G}$ represents a factual statement with subject $s \in \mathcal{E}$,
203 predicate $p \in \mathcal{R}$, and object $o \in \mathcal{E}$. For a given seed entity $e \in \mathcal{E}$, we denote the associated subgraph
204 as $\mathcal{G}_e \subseteq \mathcal{G}$, and its state after t rounds of expansion as $\mathcal{G}_e^{(t)}$.

205 4.1 Seed Subgraph Initialization.

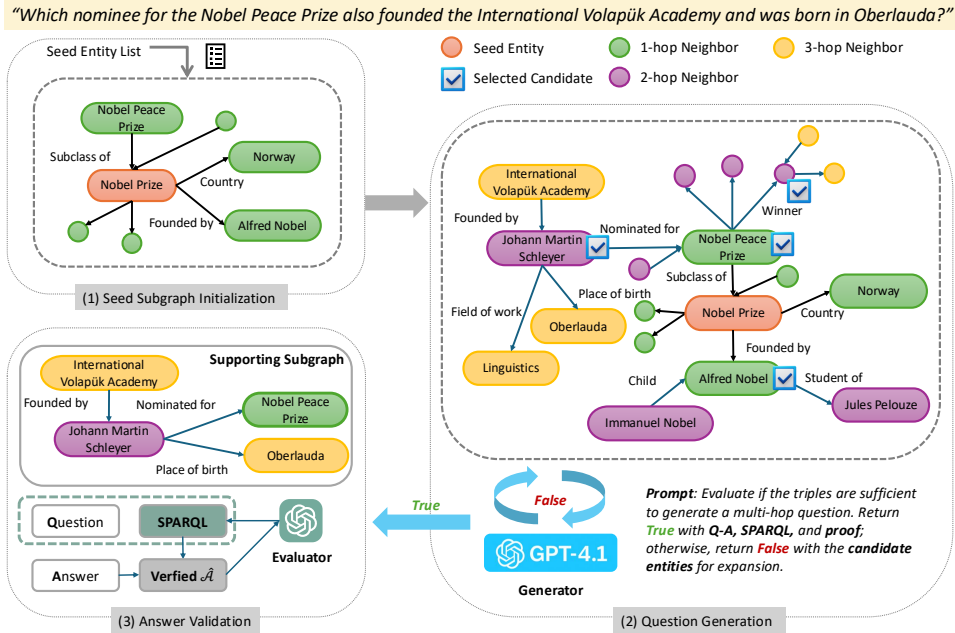
206 To ensure topic diversity and coverage across different domains, the selection of seed entities plays a
207 critical role. We draw seed entities from the Wikipedia Vital Articles³, which includes a curated set
208 of globally relevant topics spanning a broad range of domains. This provides a principled starting
209 point for constructing diverse and non-trivial question-answer pairs.

210 For each selected seed entity e , we construct an initial local subgraph $\mathcal{G}_e^{(0)}$ by sampling a fixed
211 number (e.g. 15) of its 1-hop neighbors. The resulting subgraph includes the seed and its sampled
212 neighbors, along with the triples connecting them. For example, in Figure 1, starting from the seed
213 entity *Nobel Prize*, we include sampled 1-hop neighbors such as *Nobel Peace Prize*, *Norway*, *Alfred*
214 *Nobel*, and a few other related entities, which together form the context for further expansion and
215 question generation. This initialization step helps constrain the knowledge scope while maintaining
216 enough structure to support meaningful reasoning.

217 4.2 Question Generation through Iterative LLM-Guided Subgraph Expansion

218 To generate high-quality questions that require structured, multi-hop reasoning, it is important to
219 go beyond shallow, single-fact queries. This requires subgraphs that are both semantically rich and

³https://en.wikipedia.org/wiki/Wikipedia:Vital_articles (accessed April 2, 2025)



220 structurally diverse. Therefore, we aim to expand the initialized *seed subgraph* $\mathcal{G}_e^{(0)}$ to include more
 221 informative paths and relation patterns that support deeper reasoning over the KG. A straightforward
 222 way is to expand the seed subgraph $\mathcal{G}_e^{(0)}$ using multi-hop traversal methods like BFS or DFS. However,
 223 in densely connected KGs, such unrestricted expansion quickly becomes impractical: even a few
 224 hops from a high-degree entity can yield thousands of nodes and edges, leading to bloated subgraphs
 225 that are too large to process efficiently. Moreover, these fully expanded subgraphs often contain
 226 irrelevant or weakly connected facts, making it difficult to maintain question quality.

227 Partial traversal offers a more practical alternative, typically limiting expansion to 1–2 hops or
 228 selectively sampling deep paths. But without intelligent guidance, selected paths may be arbitrary or
 229 loosely related, weakening semantic coherence. As a result, such subgraphs tend to be noisy, hard to
 230 interpret, and ill-suited for generating meaningful multi-hop questions.

231 To generate questions that are both challenging and well-formed, it is necessary to strike a balance
 232 between *coverage* (including diverse, relevant entities) and *depth* (capturing reasoning chains of
 233 sufficient complexity). To this end, we adopt an iterative expansion strategy guided by an LLM. Rather
 234 than relying on fixed-depth traversal, we allow the LLM to evaluate whether the current subgraph
 235 contains enough information to support a valid question and to suggest targeted expansion directions
 236 when additional context is needed. Once the subgraph is considered sufficient, the LLM proceeds to
 237 generate a QA instance. In the following subsections, we describe the two core components.

238 4.2.1 LLM-Guided Iterative Subgraph Expansion

239 Given an initialized subgraph $\mathcal{G}_e^{(0)}$ centered on a seed entity e , our goal is to iteratively expand it to
 240 accumulate sufficient contextual knowledge for question generation. In this subsection, we describe
 241 the $(t+1)$ -th iteration of this process, assuming that the subgraph $\mathcal{G}_e^{(t)}$ generated in the previous
 242 iteration is considered insufficient by the LLM. In this case, in the previous iteration t , the LLM also
 243 produces an *Exploration Set* $\mathcal{C}_e^{(t)}$, consisting of entities within $\mathcal{G}_e^{(t)}$ that are identified as requiring
 244 further exploration. Next, we first describe the one-hop expansion procedure for iteration $t + 1$. The
 245 LLM-based sufficiency judgment and the generation of the *Exploration Set* are discussed in detail in
 246 the Subsection of **Sufficiency Check and Exploration Set Generation**.

247 **One-Hop Expansion in Iteration $t + 1$.** Given the subgraph $\mathcal{G}_e^{(t)}$ and the *Exploration Set* $\mathcal{C}_e^{(t)}$
 248 identified by the LLM in iteration t , we expand the subgraph by incorporating additional knowledge
 249 from the global KG \mathcal{G} . Specifically, we perform a one-hop expansion around each entity $c \in \mathcal{C}_e^{(t)}$ to
 250 gather new facts that may help support the generation of more complex and informative questions.

251 For each entity c in the *Exploration Set* $\mathcal{C}_e^{(t)}$, we randomly sample a small number of its 1-hop
 252 neighbors (e.g., 10–15) and include the corresponding triples that connect these neighbors to c . This
 253 process adds a bounded amount of new information to the subgraph while keeping the expansion
 254 efficient and focused. The updated subgraph is defined as: $\mathcal{G}_e^{(t+1)} = \mathcal{G}_e^{(t)} \cup \text{SampledTriples}(\mathcal{C}_e^{(t)})$,
 255 where $\text{SampledTriples}(\mathcal{C}_e^{(t)})$ denotes the union of the selected 1-hop triples associated with the
 256 entities in $\mathcal{C}_e^{(t)}$. This one-hop expansion ensures that the subgraph grows in a controlled manner. As a
 257 result, the resulting subgraph $\mathcal{G}_e^{(t+1)}$ captures richer relational context while preserving locality and
 258 relevance to the question generation task.

259 **Sufficiency Check and Exploration Set Generation.** In this subsection, we check whether the
 260 updated subgraph $\mathcal{G}_e^{(t+1)}$ includes enough information to support a meaningful and well-scoped
 261 question. To do this, we prompt an LLM to perform two tasks: (1) evaluate whether the current
 262 subgraph is sufficient for question generation, and (2) if not, suggest a set of entities for further
 263 exploration. We refer to this set as the *Exploration Set* $\mathcal{C}_e^{(t+1)}$.

264 For *Sufficiency Checking*, we require that the subgraph support at least 2-hop reasoning and contain
 265 semantically meaningful paths sufficient to generate a well-scoped, unambiguous question. It should
 266 go beyond shallow or generic facts and provide just enough context to support a non-trivial, grounded
 267 question. The *Exploration Set* guides further expansion toward subgraph regions more likely to yield
 268 such questions. We prioritize semantically specific and structurally central entities—such as notable
 269 people, events, or awards—over broad or generic ones like countries. For example, in Figure 1, we
 270 prefer expanding on *Nobel Peace Prize* or *Alfred Nobel* over *Norway*.

271 To support both sufficiency checking and exploration set generation, we prompt GPT-4.1 with a
 272 comprehensive template. The prompt contains clear instructions and concrete examples illustrating
 273 sufficient vs. insufficient subgraphs, as well as good and bad exploration targets. The detailed prompt
 274 can be found in Appendix D.1. We will revisit this prompt in Section 4.2.2, where we describe how it
 275 also handles question and answer generation.

276 4.2.2 Question Generation from the Finalized Subgraph

277 Once a subgraph is judged to be sufficient, we denote it as \mathcal{G}_e^* , representing the finalized subgraph
 278 used for question generation. Given the finalized subgraph \mathcal{G}_e^* , we prompt the LLM to generate a
 279 complete question-answer instance. The outputs include: (1) a natural language question q_e , (2) an
 280 answer set \mathcal{A}_e , (3) a *supporting subgraph* $\mathcal{P}_e \subseteq \mathcal{G}_e^*$, and (4) a corresponding SPARQL query \mathcal{Q}_e .
 281 These components are generated jointly to maintain consistency and alignment with the reasoning
 282 required by the question. The *supporting subgraph* captures the minimal set of facts needed to justify
 283 the answer. It also facilitates error diagnosis within KGQAGen. The SPARQL query is expected to
 284 verify the answer set with the KG.

285 To ensure the generated question is meaningful and challenging, we impose several constraints. The
 286 question must be answerable using only the given finalized subgraph \mathcal{G}_e^* and involve at least two-hop
 287 reasoning. It should be specific, unambiguous, and self-contained. The question should be phrased
 288 naturally and fluently. These requirements are enforced in the unified prompt (detailed in Appendix
 289 D.1) introduced in the previous Section, which also handles sufficiency checking and exploration set
 290 selection. These three tasks are tightly linked: sufficiency checking determines whether to proceed
 291 with question generation or continue expanding the subgraph. By handling them together, we ensure
 292 that each decision is made in a shared context, leading to more coherent and aligned outputs.

293 4.3 Answer Validation and Refinement

294 The goal of this stage is to ensure that each generated question–answer pair is faithfully grounded in
 295 the KG. To achieve this, we use the SPARQL query as a formal verification tool: it must successfully
 296 retrieve the intended answer when executed over the KG.

297 Given a generated instance consisting of a question q_e , answer set \mathcal{A}_e , supporting subgraph \mathcal{P}_e , and
 298 SPARQL query \mathcal{Q}_e , we first execute \mathcal{Q}_e over the KG to retrieve a result set $\hat{\mathcal{A}}_e$. We then compare
 299 this result set to the LLM-generated answer set \mathcal{A}_e . If $\mathcal{A}_e = \hat{\mathcal{A}}_e$, we accept the instance as valid.
 300 If the two answer sets do not match, we prompt a lightweight LLM (GPT-4o-mini) to revise the
 301 SPARQL query (see Appendix D.2 for the prompts). We then re-execute the revised query and repeat

302 the validation. This revision loop continues for up to three attempts. If the revised SPARQL query
 303 returns results that match the original answer set, we retain the instance and include the revised query
 304 in the final output. Otherwise, the instance is discarded. This conservative filtering ensures that
 305 only verifiable and KG-grounded instances are retained. Additionally, since each question is paired
 306 with an executable query, the dataset can be periodically revalidated as the KG evolves—making
 307 KGQAGen-generated benchmarks both accurate at creation and maintainable over time.

308 5 KGQAGen-10k: A Sample Dataset Generated by KGQAGen

309 To demonstrate the practical capabilities of KGQAGen, we construct KGQAGen-10k, a 10K-scale
 310 KGQA dataset with KGQAGen. This dataset serves as a representative output that illustrates the
 311 effectiveness of KGQAGen in producing challenging, well-grounded, and verifiable question-answer
 312 instances. Note that the design of KGQAGen is highly modular and scalable. With minimal human
 313 effort, it can be readily extended to generate substantially larger datasets across diverse knowledge
 314 domains. In this section, we focus on the construction and analysis of KGQAGen-10k.

315 **Dataset Construction.** To construct KGQAGen-10k, we begin by selecting 16,000 seed entities from
 316 Wikipedia’s Level-5 Vital Articles, which provide broad topical coverage and global relevance. All
 317 questions are grounded in the most up-to-date dump of Wikidata⁴, a large and publicly accessible
 318 knowledge base with comprehensive entity and relation coverage. Using KGQAGen, we generate
 319 15,451 question-answer pairs. After applying answer validation to remove examples with non-
 320 executable or inconsistent SPARQL queries, we obtain 10,787 verified instances. We refer to this
 321 dataset as KGQAGen-10k.

322 **Manual Quality Assessment.** To assess the factual accuracy of
 323 KGQAGen-10k, we conducted a manual audit of 300 randomly
 324 sampled question-answer pairs. Following the same annotation
 325 protocol used in Section 3, each example was carefully
 326 verified. We found that 289 out of 300 examples (96.3%) are
 327 correct. This result suggests that the QA instances generated
 328 by KGQAGen are highly reliable and well-grounded in verifi-
 329 able knowledge. We defer analysis of question difficulty to
 330 Section 6, where we benchmark several LLM and KG-RAG
 331 models. A case study of the generated questions is provided in
 332 Appendix F.2.

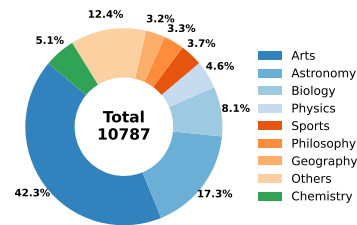


Figure 2: Topic Distribution

333 **Statistics and Properties.** We analyze KGQAGen-10k along two key dimensions: linguistic complex-
 334 ity and topic coverage: (1) *Linguistic Complexity*: Most questions (61.1%) are 16–30 words long,
 335 showing moderate complexity. A third exceed 30 words, indicating deeper reasoning. Only 7.5%
 336 are short factoid queries, highlighting the dataset’s focus on rich, natural questions. For answer set
 337 size, 84.5% of the questions yield a single answer, while 9.7% return three or more answers. (2)
 338 *Structural Difficulty*: Following DyVal [81], we compute six graph-based metrics for each supporting
 339 subgraph—number of nodes, edges, average degree, depth (hops), width (maximum frontier), and
 340 extra links (non-chain shortcuts). Across all 10,787 examples, 98% require 2–5 hops, 84% contain
 341 5–30 entities, 83% have 4–28 relations, 79% exhibit width 4–20, and 92% are fully connected. These
 342 statistics confirm that KGQAGen-10k consists of diverse, non-trivial reasoning graphs with substantial
 343 structural complexity. A full summary table is provided in Appendix F.1. (3) *Topic Coverage*: The
 344 dataset also provides broad semantic coverage. Figure 2 shows that KGQAGen-10k covers a broad
 345 range of topics. The most represented categories are Arts (42.3%) and Astronomy (17.3%), followed
 346 by STEM fields (16% combined), Sports, Geography, and Philosophy.

347 6 KGQAGen-10k as a Benchmark for KG-RAG Models

348 In this section, we use KGQAGen-10k to benchmark a variety of KG-RAG models and LLMs.

349 **Experiment Setup.** We benchmark all models on KGQAGen-10k using a standardized split of 8,629
 350 training, 1,079 development, and 1,079 test examples. Most prior work evaluates KGQA systems
 351 using an exact string match between predicted and reference answers. However, as highlighted in

⁴<https://dumps.wikimedia.org/wikidatawiki/> (accessed April 2, 2025)

Table 2: Performance on KGQAGen-10k. EM = Exact Match; LASM = LLM-Assisted Semantic Match; LLM-SP = LLM with Supporting Subgraph as input.

Type	Model	Accuracy		Hit@1		F1		Precision		Recall	
		EM	LASM	EM	LASM	EM	LASM	EM	LASM	EM	LASM
Pure LLM	LLaMA-3.1-8B-Instruct [19]	8.87	11.91	9.27	12.42	8.97	11.98	9.27	12.42	8.87	11.81
	LLaMA2-7B [61]	6.88	12.32	7.23	12.88	6.95	12.34	7.23	13.72	6.88	11.96
	Mistral-7B-Instruct-v0.2 [2]	24.98	32.34	26.51	34.38	25.36	33.20	26.60	34.85	24.98	32.72
	GPT-4o-mini [1]	32.34	42.49	34.11	44.39	32.74	42.91	34.11	44.86	32.34	42.35
	GPT-4 [1]	42.38	51.37	44.95	54.49	43.01	52.32	45.13	55.33	42.38	51.48
	DeepSeek-Chat [11]	42.48	51.84	45.51	55.24	43.17	52.64	45.60	55.79	42.48	51.78
	GPT-4o [1]	45.29	54.21	47.91	57.46	45.89	54.93	48.01	57.83	45.29	54.11
	GPT-4.1 [41]	47.43	56.96	50.05	59.96	48.03	57.72	50.05	60.33	47.43	56.95
RAG-based	RoG (LLaMA2-7B) [33]	20.10	27.28	21.32	28.92	17.75	24.26	17.65	24.79	20.10	27.16
	GCR (LLaMA-3.1 + GPT-4o) [34]	49.37	58.96	52.46	62.84	49.30	58.88	50.61	60.76	49.37	59.18
	ToG (GPT-4o) [57]	49.65	59.89	52.55	63.02	50.38	60.73	53.01	64.23	49.65	59.76
	PoG (GPT-4o) [9]	50.67	60.18	54.03	63.95	51.47	61.30	54.31	64.78	50.67	60.34
LLM-SP	LLaMA2-7B (w/ SP) [61]	69.79	73.79	73.12	77.76	70.43	74.55	72.81	77.67	69.79	73.76
	GPT-4o (w/ SP) [1]	82.46	84.89	89.62	92.22	84.07	86.75	89.81	93.23	82.46	84.95

352 Section 3, this metric often fails to capture semantically valid predictions that differ in surface form.
353 To address this limitation, we introduce **LLM-Assisted Semantic Match (LASM)**, a lightweight
354 verification mechanism that activates only when a prediction fails the exact match. LASM uses
355 GPT-4o-mini to assess semantic equivalence between predictions and gold answers (details in
356 Appendix G). We report Accuracy, Hit@1, Precision, Recall, and F1 under both Exact Match (EM)
357 and LASM for direct comparison.

358 **Evaluated Models.** We evaluate three categories of models on KGQAGen-10k: (1) *Pure LLMs*:
359 This group includes a range of open-source and commercial models—LLaMA-3.1-8B-Instruct,
360 LLaMA2-7B, Mistral-7B-Instruct-v0.2, GPT-4o-mini, GPT-4, GPT-4o, GPT-4.1, and
361 Deepseek-Chat. These models receive only the natural language question without any retrieval or
362 external grounding; (2) *KG-RAG Models*: We include recent KG-RAG models such as RoG, GCR, ToG
363 , and PoG. These models use the KG as an external retrieval source, either by incorporating retrieved
364 triples, augmenting the prompt with symbolic paths, or integrating topology-aware context; and (3)
365 *LLM with Supporting Subgraph (LLM-SP)*. To assess performance under perfect retrieval, we include
366 LLaMA2-7B and GPT-4o models that are directly given the associated *supporting subgraph* of the
367 questions from the dataset construction process. These setups simulate perfect retrieval and help
368 isolate the reasoning capabilities from retrieval effectiveness. See Appendix G for more details of
369 baseline models.

370 **Results and Analysis.** Table 2 summarizes the performance of all evaluated models across standard
371 metrics under both EM and LASM protocols. We make the following key observations: (1) Even
372 capable LLMs such as GPT-4.1 and recent KG-RAG models like GCR and PoG achieve only moderate
373 performance on KGQAGen-10k, showing the non-trivial nature of the questions and the challenges
374 they pose for current methods. (2) LASM consistently yields higher reported performance across all
375 models over EM, indicating that many predictions marked incorrect under exact match are actually
376 semantically correct. This highlights the importance of incorporating semantic-aware evaluation
377 to more accurately assess QA performance. (3) Model capability strongly correlates with perfor-
378 mance: larger and more recent LLMs such as GPT-4.1 substantially outperform smaller models
379 like LLaMA2-7B, showing the importance of both knowledge coverage and reasoning ability. (4)
380 KG-RAG models achieve noticeable gains over their corresponding LLM backbones. For example,
381 RoG improves upon LLaMA2-7B by over 10 points, while GCR, ToG, and PoG outperform GPT-4o
382 by approximately 4 points. These results confirm that incorporating external KG-derived context
383 enhances QA performance. However, the overall improvement remains moderate, suggesting that
384 retrieval components in current KG-RAG systems are still suboptimal and warrant further enhance-
385 ment. (5) LLM-SP models—LLMs provided with the ground truth *supporting subgraph*—achieve
386 the strongest performance by a substantial margin. For instance, LLaMA2-7B (w/ SP) reaches LASM
387 accuracy 73.8%, outperforming not only its base model LLaMA2-7B but all other larger LLMs. Simi-
388 larly, GPT-4o (w/ SP) sees its LASM accuracy rise from 54.2% to 84.9% compared with GPT-4o.
389 These results highlight the key role of high-quality retrieval in KG-RAG systems and suggest that
390 retrieval remains a major bottleneck limiting current KG-RAG systems.

391 **Cross-Dataset Comparison.** To further contextualize the results, we evaluate representative KG-
392 RAG models on two established KGQA benchmarks, WebQSP [73] and CWQ [58], in addition to

393 our KGQAGen-10k. While models such as RoG, GCR, ToG, and PoG achieve high Hit@1 scores on We-
394 bQSP and CWQ (e.g., 85.7 and 92.2, respectively), their performance drops sharply on KGQAGen-10k
395 (21.3–54.0), reflecting its higher reasoning complexity and stricter grounding requirements. De-
396 tailed results are provided in Appendix G.3. This consistent decline across models highlights that
397 legacy benchmarks often overestimate true reasoning ability, whereas KGQAGen-10k offers a more
398 challenging and diagnostic testbed for evaluating KG-RAG systems.

399 7 Conclusion

400 In this paper, we identified major quality issues in existing KGQA benchmarks through a detailed
401 manual audit, including inaccurate answers and poorly constructed questions. To address this,
402 we introduced KGQAGen, a framework for building well-grounded, verifiable QA datasets at scale.
403 Using this framework, we created KGQAGen-10k, a 10K-example benchmark that challenges both
404 general-purpose LLMs and KG-RAG models. Our results show that even strong models struggle
405 on KGQAGen-10k, highlighting the need for better retrieval and reasoning mechanisms. They also
406 demonstrate the value of KGQAGen as a scalable and effective approach for constructing challenging,
407 high-quality benchmarks to support future progress for KG-RAG systems.

408 **Limitations:** While KGQAGen enables scalable and verifiable KGQA dataset construction, it relies
409 on the availability and accuracy of the underlying KG (e.g., Wikidata). Errors or gaps in the KG
410 can limit the quality of generated questions. In addition, our pipeline depends on LLM capabilities
411 and prompt design; although we mitigate hallucination risks via symbolic verification, the quality of
412 subgraph selection and question phrasing may still be influenced by LLM variability.

413 Acknowledgments

414 This research is supported by the National Science Foundation (NSF) under grant numbers NSF-
415 2406647 and NSF-2406648. It is also supported by the National Artificial Intelligence Research
416 Resource (NAIRR) Pilot and the Delta advanced computing and data resource, which is supported by
417 the National Science Foundation under award NSF-OAC-2005572.

418 References

- 419 [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
420 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4
421 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 422 [2] Mistral AI. Models overview, 2025. Accessed: 11 May 2025.
- 423 [3] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary
424 Ives. Dbpedia: A nucleus for a web of open data. In *international semantic web conference*,
425 pages 722–735. Springer, 2007.
- 426 [4] Michael Azmy, Peng Shi, Jimmy Lin, and Ihab Ilyas. Farewell freebase: Migrating the
427 simplequestions dataset to dbpedia. In *Proceedings of the 27th international conference on*
428 *computational linguistics*, pages 2093–2103, 2018.
- 429 [5] Junwei Bao, Nan Duan, Zhao Yan, Ming Zhou, and Tiejun Zhao. Constraint-based question
430 answering with knowledge graph. In *Proceedings of COLING 2016, the 26th international*
431 *conference on computational linguistics: technical papers*, pages 2503–2514, 2016.
- 432 [6] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase
433 from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in*
434 *natural language processing*, pages 1533–1544, 2013.
- 435 [7] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a
436 collaboratively created graph database for structuring human knowledge. In *Proceedings of the*
437 *2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250, 2008.

- 438 [8] Shulin Cao, Jiaxin Shi, Liangming Pan, Lunyiu Nie, Yutong Xiang, Lei Hou, Juanzi Li, Bin He,
439 and Hanwang Zhang. Kqa pro: A dataset with explicit compositional programs for complex
440 question answering over knowledge base. *arXiv preprint arXiv:2007.03875*, 2020.
- 441 [9] Liyi Chen, Panrong Tong, Zhongming Jin, Ying Sun, Jieping Ye, and Hui Xiong. Plan-on-graph:
442 Self-correcting adaptive planning of large language model on knowledge graphs. *arXiv preprint*
443 *arXiv:2410.23875*, 2024.
- 444 [10] Preetam Prabhu Srikar Dammu, Himanshu Naidu, and Chirag Shah. Dynamic-kgqa: A
445 scalable framework for generating adaptive question answering datasets. *arXiv preprint*
446 *arXiv:2503.05049*, 2025.
- 447 [11] DeepSeek-AI. Deepseek-v3 technical report, 2024.
- 448 [12] Mohammad Dehghan, Mohammad Ali Alomrani, Sunyam Bagga, David Alfonso-Hermelo,
449 Khalil Bibi, Abbas Ghaddar, Yingxue Zhang, Xiaoguang Li, Jianye Hao, Qun Liu, et al. Ewek-
450 qa: Enhanced web and efficient knowledge graph retrieval for citation-based question answering
451 systems. *arXiv preprint arXiv:2406.10393*, 2024.
- 452 [13] Zixuan Dong, Baoyun Peng, Yufei Wang, Jia Fu, Xiaodong Wang, Yongxue Shan, and Xin Zhou.
453 Effiqa: Efficient question-answering with strategic multi-model collaboration on knowledge
454 graphs. *arXiv preprint arXiv:2406.01238*, 2024.
- 455 [14] Mohnish Dubey, Debayan Banerjee, Abdelrahman Abdelkawi, and Jens Lehmann. Lc-quad
456 2.0: A large dataset for complex question answering over wikidata and dbpedia. In *The*
457 *Semantic Web–ISWC 2019: 18th International Semantic Web Conference, Auckland, New*
458 *Zealand, October 26–30, 2019, Proceedings, Part II 18*, pages 69–78. Springer, 2019.
- 459 [15] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven
460 Truitt, Dasha Metropolitanaky, Robert Osazuwa Ness, and Jonathan Larson. From local to global:
461 A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- 462 [16] Haishuo Fang, Xiaodan Zhu, and Iryna Gurevych. Dara: Decomposition-alignment-reasoning
463 autonomous language agent for question answering over knowledge graphs. *arXiv preprint*
464 *arXiv:2406.07080*, 2024.
- 465 [17] Siyuan Fang, Kaijing Ma, Tianyu Zheng, Xinrun Du, Ningxuan Lu, Ge Zhang, and Qingkun
466 Tang. Karpa: A training-free method of adapting knowledge graph as references for large
467 language model’s reasoning path aggregation. *arXiv preprint arXiv:2412.20995*, 2024.
- 468 [18] Zengyi Gao, Yukun Cao, Hairu Wang, Ao Ke, Yuan Feng, Xike Xie, and S Kevin Zhou. Frag:
469 A flexible modular framework for retrieval-augmented generation based on knowledge graphs.
470 *arXiv preprint arXiv:2501.09957*, 2025.
- 471 [19] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian,
472 Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama
473 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 474 [20] Yu Gu, Sue Kase, Michelle Vanni, Brian Sadler, Percy Liang, Xifeng Yan, and Yu Su. Beyond
475 iid: three levels of generalization for question answering on knowledge bases. In *Proceedings*
476 *of the Web Conference 2021*, pages 3477–3488, 2021.
- 477 [21] Gaole He, Yunshi Lan, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. Improving multi-
478 hop knowledge base question answering by learning intermediate supervision signals. In
479 *Proceedings of the 14th ACM international conference on web search and data mining*, pages
480 553–561, 2021.
- 481 [22] Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. Grag: Graph
482 retrieval-augmented generation. *arXiv preprint arXiv:2405.16506*, 2024.
- 483 [23] Kelvin Jiang, Dekun Wu, and Hui Jiang. Freebaseqa: A new factoid qa data set matching
484 trivia-style question-answer pairs with freebase. In *Proceedings of the 2019 Conference of the*
485 *North American Chapter of the Association for Computational Linguistics: Human Language*
486 *Technologies, Volume 1 (Long and Short Papers)*, pages 318–323, 2019.

- 487 [24] Longquan Jiang and Ricardo Usbeck. Knowledge graph question answering datasets and their
488 generalizability: Are they enough for future research? In *Proceedings of the 45th International*
489 *ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3209–
490 3218, 2022.
- 491 [25] Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashu-
492 bin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, et al. Measuring
493 compositional generalization: A comprehensive method on realistic data. *arXiv preprint*
494 *arXiv:1912.09713*, 2019.
- 495 [26] Kun Li, Tianhua Zhang, Xixin Wu, Hongyin Luo, James Glass, and Helen Meng. Decoding on
496 graphs: Faithful and sound reasoning on knowledge graphs through generation of well-formed
497 chains. *arXiv preprint arXiv:2410.18415*, 2024.
- 498 [27] Mufei Li, Siqi Miao, and Pan Li. Simple is effective: The roles of graphs and large language mod-
499 els in knowledge-graph-based retrieval-augmented generation. *arXiv preprint arXiv:2410.20724*,
500 2024.
- 501 [28] Biyang Liu, Huimin Yu, and Guodong Qi. Graftnet: Towards domain generalized stereo
502 matching with a broad-spectrum and task-oriented feature. In *Proceedings of the IEEE/CVF*
503 *conference on computer vision and pattern recognition*, pages 13012–13021, 2022.
- 504 [29] Guangyi Liu, Yongqi Zhang, Yong Li, and Quanming Yao. Dual reasoning: A gnn-llm collabora-
505 tive framework for knowledge graph question answering. *Proceedings of the Conference on*
506 *Parsimony and Learning (CPAL)*, 2024.
- 507 [30] Guangyi Liu, Yongqi Zhang, Yong Li, and Quanming Yao. Dual reasoning: A gnn-llm
508 collaborative framework for knowledge graph question answering. In *The Second Conference*
509 *on Parsimony and Learning (Proceedings Track)*, 2024.
- 510 [31] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval:
511 Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*,
512 2023.
- 513 [32] Haoran Luo, Zichen Tang, Shiyao Peng, Yikai Guo, Wentai Zhang, Chenghao Ma, Guanting
514 Dong, Meina Song, Wei Lin, Yifan Zhu, et al. Chatkbqa: A generate-then-retrieve framework
515 for knowledge base question answering with fine-tuned large language models. *arXiv preprint*
516 *arXiv:2310.08975*, 2023.
- 517 [33] Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. Reasoning on graphs: Faithful
518 and interpretable large language model reasoning. *arXiv preprint arXiv:2310.01061*, 2023.
- 519 [34] Linhao Luo, Zicheng Zhao, Chen Gong, Gholamreza Haffari, and Shirui Pan. Graph-constrained
520 reasoning: Faithful reasoning on knowledge graphs with large language models. *arXiv preprint*
521 *arXiv:2410.13080*, 2024.
- 522 [35] Costas Mavromatis and George Karypis. Rearev: Adaptive reasoning for question answering
523 over knowledge graphs. *arXiv preprint arXiv:2210.13650*, 2022.
- 524 [36] Costas Mavromatis and George Karypis. Gnn-rag: Graph neural retrieval for large language
525 model reasoning. *arXiv preprint arXiv:2405.20139*, 2024.
- 526 [37] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit
527 Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation
528 of factual precision in long form text generation. *arXiv preprint arXiv:2305.14251*, 2023.
- 529 [38] John Paul Mueller. *Mining Google web services: building applications with the Google API*.
530 John Wiley & Sons, 2006.
- 531 [39] Ngonga Ngomo. 9th challenge on question answering over linked data (qald-9). *language*,
532 7(1):58–64, 2018.
- 533 [40] Reham Omar, Omij Mangukiya, and Essam Mansour. Dialogue benchmark generation from
534 knowledge graphs with cost-effective retrieval-augmented llms. *Proceedings of the ACM on*
535 *Management of Data*, 3(1):1–26, 2025.

- 536 [41] OpenAI. Introducing gpt-4.1 in the api, April 2025. Accessed: 2025-05-11.
- 537 [42] Abdelghny Orogat and Ahmed El-Roby. Smartbench: demonstrating automatic generation of
538 comprehensive benchmarks for question answering over knowledge graphs. *Proceedings of the*
539 *VLDB Endowment*, 15(12):3662–3665, 2022.
- 540 [43] Abdelghny Orogat and Ahmed El-Roby. Maestro: Automatic generation of comprehensive
541 benchmarks for question answering over knowledge graphs. *Proceedings of the ACM on*
542 *Management of Data*, 1(2):1–24, 2023.
- 543 [44] Abdelghny Orogat, Isabelle Liu, and Ahmed El-Roby. Cbench: Towards better evaluation of
544 question answering over knowledge graphs. *arXiv preprint arXiv:2105.00811*, 2021.
- 545 [45] Liangming Pan, Wenhui Chen, Min-Yen Kan, and William Yang Wang. Attacking open-domain
546 question answering by injecting misinformation. *arXiv preprint arXiv:2110.07803*, 2021.
- 547 [46] Thomas Pellissier Tanon, Denny Vrandečić, Sebastian Schaffert, Thomas Steiner, and Lydia
548 Pintscher. From freebase to wikidata: The great migration. In *Proceedings of the 25th*
549 *international conference on world wide web*, pages 1419–1428, 2016.
- 550 [47] Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang,
551 and Siliang Tang. Graph retrieval-augmented generation: A survey. *arXiv preprint*
552 *arXiv:2408.08921*, 2024.
- 553 [48] Aleksandr Perevalov, Dennis Diefenbach, Ricardo Usbeck, and Andreas Both. Qald-9-plus:
554 A multilingual dataset for question answering over dbpedia and wikidata translated by native
555 speakers. In *2022 IEEE 16th International Conference on Semantic Computing (ICSC)*, pages
556 229–234. IEEE, 2022.
- 557 [49] Tyler Thomas Procko and Omar Ochoa. Graph retrieval-augmented generation for large
558 language models: A survey. In *2024 Conference on AI, Science, Engineering, and Technology*
559 *(AIxSET)*, pages 166–169. IEEE, 2024.
- 560 [50] Julian Risch, Timo Möller, Julian Gutsch, and Malte Pietsch. Semantic answer similarity for
561 evaluating question answering models. *arXiv preprint arXiv:2108.06130*, 2021.
- 562 [51] Amrita Saha, Vardaan Pahuja, Mitesh Khapra, Karthik Sankaranarayanan, and Sarath Chandar.
563 Complex sequential question answering: Towards learning to converse over linked question
564 answer pairs with a knowledge graph. In *Proceedings of the AAAI conference on artificial*
565 *intelligence*, volume 32, 2018.
- 566 [52] Tiesunlong Shen, Jin Wang, Xuejie Zhang, and Erik Cambria. Reasoning with trees: Faithful
567 question answering over knowledge graph. In *Proceedings of the 31st International Conference*
568 *on Computational Linguistics*, pages 3138–3157, 2025.
- 569 [53] Manthankumar Solanki. Efficient document retrieval with g-retriever. *arXiv preprint*
570 *arXiv:2504.14955*, 2025.
- 571 [54] Yu Su, Huan Sun, Brian Sadler, Mudhakar Srivatsa, Izzeddin Gür, Zenghui Yan, and Xifeng
572 Yan. On generating characteristic-rich question sets for qa evaluation. In *Proceedings of the*
573 *2016 Conference on Empirical Methods in Natural Language Processing*, pages 562–572, 2016.
- 574 [55] Fabian M Suchanek, Mehwish Alam, Thomas Bonald, Lihu Chen, Pierre-Henri Paris, and Jules
575 Soria. Yago 4.5: A large and clean knowledge base with a rich taxonomy. In *Proceedings of*
576 *the 47th International ACM SIGIR Conference on Research and Development in Information*
577 *Retrieval*, pages 131–140, 2024.
- 578 [56] Haitian Sun, Tania Bedrax-Weiss, and William W Cohen. Pullnet: Open domain question
579 answering with iterative retrieval on knowledge bases and text. *arXiv preprint arXiv:1904.09537*,
580 2019.
- 581 [57] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel M
582 Ni, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of
583 large language model on knowledge graph. *arXiv preprint arXiv:2307.07697*, 2023.

- 584 [58] Alon Talmor and Jonathan Berant. The web as a knowledge-base for answering complex
585 questions. *arXiv preprint arXiv:1803.06643*, 2018.
- 586 [59] Xingyu Tan, Xiaoyang Wang, Qing Liu, Xiwei Xu, Xin Yuan, and Wenjie Zhang. Paths-over-
587 graph: Knowledge graph empowered large language model reasoning. In *Proceedings of the*
588 *ACM on Web Conference 2025*, pages 3505–3522, 2025.
- 589 [60] Xiaqiang Tang, Jian Li, Nan Du, and Sihong Xie. Adapting to non-stationary environments:
590 Multi-armed bandit enhanced retrieval-augmented generation on knowledge graphs. In *Pro-*
591 *ceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 12658–12666,
592 2025.
- 593 [61] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
594 Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
595 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 596 [62] Priyansh Trivedi, Gaurav Maheshwari, Mohnish Dubey, and Jens Lehmann. Lc-quad: A
597 corpus for complex question answering over knowledge graphs. In *The Semantic Web–ISWC*
598 *2017: 16th International Semantic Web Conference, Vienna, Austria, October 21–25, 2017,*
599 *Proceedings, Part II 16*, pages 210–218. Springer, 2017.
- 600 [63] Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Com-*
601 *munications of the ACM*, 57(10):78–85, 2014.
- 602 [64] Cunxiang Wang, Sirui Cheng, Qipeng Guo, Yuanhao Yue, Bowen Ding, Zhikun Xu, Yidong
603 Wang, Xiangkun Hu, Zheng Zhang, and Yue Zhang. Evaluating open-qa evaluation. *Advances*
604 *in Neural Information Processing Systems*, 36:77013–77042, 2023.
- 605 [65] Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge
606 Rong, and Zhang Xiong. Knowledge-driven cot: Exploring faithful reasoning in llms for
607 knowledge-intensive question answering. *arXiv preprint arXiv:2308.13259*, 2023.
- 608 [66] Song Wang, Junhong Lin, Xiaojie Guo, Julian Shun, Jundong Li, and Yada Zhu. Reason-
609 ing of large language models over knowledge graphs with super-relations. *arXiv preprint*
610 *arXiv:2503.22166*, 2025.
- 611 [67] Liqiang Wen, Guanming Xiong, Tong Mo, Bing Li, Weiping Li, and Wen Zhao. Clear-kgqa:
612 Clarification-enhanced ambiguity resolution for knowledge graph question answering. *arXiv*
613 *preprint arXiv:2504.09665*, 2025.
- 614 [68] Alexander R Fabbri Chien-Sheng Wu and Wenhao Liu Caiming Xiong. Qafacteval: Improved
615 qa-based factual consistency evaluation for summarization. 2023.
- 616 [69] Guanming Xiong, Junwei Bao, and Wen Zhao. Interactive-kbqa: Multi-turn interac-
617 tions for knowledge base question answering with large language models. *arXiv preprint*
618 *arXiv:2402.15131*, 2024.
- 619 [70] Mufan Xu, Kehai Chen, Xuefeng Bai, Muyun Yang, Tiejun Zhao, and Min Zhang. Llm-
620 based discriminative reasoning for knowledge graph question answering. *arXiv preprint*
621 *arXiv:2412.12643*, 2024.
- 622 [71] Mufan Xu, Gewen Liang, Kehai Chen, Wei Wang, Xun Zhou, Muyun Yang, Tiejun Zhao,
623 and Min Zhang. Memory-augmented query reconstruction for llm-based knowledge graph
624 reasoning. *arXiv preprint arXiv:2503.05193*, 2025.
- 625 [72] Yao Xu, Shizhu He, Jiabei Chen, Zihao Wang, Yangqiu Song, Hanghang Tong, Guang Liu,
626 Kang Liu, and Jun Zhao. Generate-on-graph: Treat llm as both agent and kg in incomplete
627 knowledge graph question answering. *arXiv preprint arXiv:2404.14741*, 2024.
- 628 [73] Wen-tau Yih, Matthew Richardson, Christopher Meek, Ming-Wei Chang, and Jina Suh. The
629 value of semantic parse labeling for knowledge base question answering. In *Proceedings of*
630 *the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short*
631 *Papers)*, pages 201–206, 2016.

- 632 [74] Donghan Yu, Sheng Zhang, Patrick Ng, Henghui Zhu, Alexander Hanbo Li, Jun Wang, Yiqun
633 Hu, William Wang, Zhiguo Wang, and Bing Xiang. Decaf: Joint decoding of answers and
634 logical forms for question answering over knowledge bases. *arXiv preprint arXiv:2210.00063*,
635 2022.
- 636 [75] Buchao Zhan, Anqi Li, Xin Yang, Dongmei He, Yucong Duan, and Shankai Yan. Rarok:
637 Retrieval-augmented reasoning on knowledge for medical question answering. In *2024 IEEE*
638 *International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 2837–2843. IEEE,
639 2024.
- 640 [76] Jing Zhang, Xiaokang Zhang, Jifan Yu, Jian Tang, Jie Tang, Cuiping Li, and Hong Chen.
641 Subgraph retrieval enhanced model for multi-hop knowledge base question answering. *arXiv*
642 *preprint arXiv:2202.13296*, 2022.
- 643 [77] Qinggang Zhang, Shengyuan Chen, Yuanchen Bei, Zheng Yuan, Huachi Zhou, Zijin Hong,
644 Junnan Dong, Hao Chen, Yi Chang, and Xiao Huang. A survey of graph retrieval-augmented
645 generation for customized large language models. *arXiv preprint arXiv:2501.13958*, 2025.
- 646 [78] Wen Zhang, Long Jin, Yushan Zhu, Jiaoyan Chen, Zhiwei Huang, Junjie Wang, Yin Hua, Lei
647 Liang, and Huajun Chen. Trustuqa: A trustful framework for unified structured data question
648 answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages
649 25931–25939, 2025.
- 650 [79] Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander Smola, and Le Song. Variational
651 reasoning for question answering with knowledge graph. In *Proceedings of the AAAI conference*
652 *on artificial intelligence*, volume 32, 2018.
- 653 [80] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
654 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
655 chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.
- 656 [81] Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie.
657 Dyval: Dynamic evaluation of large language models for reasoning tasks. *arXiv preprint*
658 *arXiv:2309.17167*, 2023.

659 **NeurIPS Paper Checklist**

660 The checklist is designed to encourage best practices for responsible machine learning research,
661 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
662 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
663 follow the references and follow the (optional) supplemental material. The checklist does NOT count
664 towards the page limit.

665 Please read the checklist guidelines carefully for information on how to answer these questions. For
666 each question in the checklist:

- 667 • You should answer [Yes] , [No] , or [NA] .
- 668 • [NA] means either that the question is Not Applicable for that particular paper or the
669 relevant information is Not Available.
- 670 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

671 **The checklist answers are an integral part of your paper submission.** They are visible to the
672 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
673 (after eventual revisions) with the final version of your paper, and its final version will be published
674 with the paper.

675 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
676 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
677 proper justification is given (e.g., "error bars are not reported because it would be too computationally
678 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
679 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
680 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
681 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
682 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
683 please point to the section(s) where related material for the question can be found.

684 IMPORTANT, please:

- 685 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 686 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 687 • **Do not modify the questions and only use the provided macros for your answers.**

688 **1. Claims**

689 Answer: [Yes]

690 Justification: The abstract and introduction clearly state our contributions, including the
691 current benchmark pitfalls analysis, our KGQAGen framework, the sampled KGQAGen-10k
692 dataset, and a systematic evaluation of baseline KGQA models. See Section 3, 4, 5, 6.

693 **2. Limitations**

694 Answer: [Yes]

695 Justification: We discuss limitations, including incomplete answer coverage due to KG gaps,
696 dependence on the accuracy and currency of the underlying KG, and sensitivity to LLM
697 behavior during graph expansion and question generation. While symbolic verification
698 helps mitigate hallucinations, the quality of subgraph selection and phrasing may still be
699 influenced by LLM variability. See Section 7: Conclusion.

700 **3. Theory assumptions and proofs**

701 Answer: [NA]

702 Justification: The paper does not include formal theoretical results or proofs, as our work is
703 empirical and dataset-centric.

704 **4. Experimental result reproducibility**

705 Answer: [Yes]

706 Justification: All implementation details, model configurations, and data processing steps
707 are provided in the supplementary material. The baseline models reproduce code is publicly
708 available at <https://github.com/liangliang6v6/KGQAGen>.

709 5. Open access to data and code

710 Answer: [Yes]

711 Justification: The KGQAGen benchmark code is publicly available at <https://github.com/liangliang6v6/KGQAGen> and the KGQAGen-10k dataset is hosted at <https://huggingface.co/datasets/lianglz/KGQAGen-10k>. Instructions are included for full
712 reproducibility.
713
714

715 6. Experimental setting/details

716 Answer: [Yes]

717 Justification: We provide comprehensive implementation details for dataset construction,
718 model evaluation in our code, including fine-tuning configurations and hyperparameters for
719 all baseline models. See Appendix G for the full experimental setup of benchmarking.

720 7. Experiment statistical significance

721 Answer: [No]

722 Justification: The reported results are based on single-run evaluations without repeated trials
723 due to the computational cost of running large-scale LLM experiments. As such, we do not
724 report statistical significance tests or error bars.

725 8. Experiments compute resources

726 Answer: [Yes]

727 Justification: While we did not rely on heavy GPU usage for dataset generation, our
728 framework involves computational costs from LLM API calls (e.g., GPT-4o) and deploying
729 a local Virtuoso endpoint over large Wiki dumps. The main GPU resources were used for
730 fine-tuning baseline models such as RoG and GCR. Details are provided in Appendix G.

731 9. Code of ethics

732 Answer: [Yes]

733 Justification: We have reviewed and followed the NeurIPS Code of Ethics. Our dataset is
734 filtered to exclude harmful content, and the methodology avoids any deceptive practices.

735 10. Broader impacts

736 Answer: [Yes]

737 Justification: Our KGQAGen framework offers a flexible pipeline for KG-RAG model
738 development and can be readily adapted to domain-specific applications. For example, it can
739 be used to generate high-quality QA datasets in fields like biology or medicine by leveraging
740 domain-specific knowledge graphs. We discuss positive impacts, e.g., improved semantic
741 evaluation, better KGQA benchmarks in Sections 5, 6.

742 Given the scope and focus of our work, we believe the potential for negative societal impacts
743 is minimal. Our framework is designed for benchmark dataset generation to support research
744 in KGQA, and is not directly tied to deployment in sensitive or high-risk applications. We
745 do not foresee direct risks related to disinformation, surveillance, or privacy, as our system
746 operates on publicly available and curated knowledge bases (e.g., Wikidata) and focuses
747 on structured question generation and evaluation. Additionally, the symbolic verification
748 and grounding aspects of our approach are designed to reduce hallucinations and improve
749 factual consistency, which could mitigate misuse in downstream applications. Nonetheless,
750 we acknowledge that if applied without proper care in downstream domains, there could be
751 potential issues due to biases or inaccuracies in the underlying knowledge graph. These risks
752 are mitigated in our framework through symbolic execution and validation mechanisms,
753 which help filter incorrect outputs.

754 11. Safeguards

755 Answer: [Yes]

756 Justification: GPT API calls and online Wikidata SPARQL requests are rate-limited to
757 prevent abuse. We use SPARQL-based answer verification to filter incorrect generations. The
758 framework is designed for research use with safeguards to ensure responsible deployment.
759

12. **Licenses for existing assets**

760 Answer: [Yes]

761 Justification: We use Wikidata (CC0 license) and OpenAI models, both properly cited with
762 license information in the supplemental material. Our KGQAGen-10k dataset is openly
763 licensed on Hugging Face, and the source code is released on GitHub under the MIT
764 license. The work does not involve crowdsourcing or human subjects, so IRB approval is
765 not applicable.

13. **New assets**

766 Answer: [Yes]

767 Justification: The KGQAGen-10k dataset is newly introduced and released under an open
768 license. It includes documentation, data format specifications, and usage instructions with
769 example, hosted on Hugging Face.
770

14. **Crowdsourcing and research with human subjects**

771 Answer: [NA]

772 Justification: No human subjects or crowd annotation were involved in the creation or
773 validation of KGQAGen-10k.
774

15. **Institutional review board (IRB) approvals or equivalent for research with human
775 subjects**

776 Answer: [NA]

777 Justification: No IRB approval is required since the research does not involve human
778 subjects.
779

16. **Declaration of LLM usage**

780 Answer: [Yes]

781 Justification: We use GPT-4 . 1 and GPT-4o-mini for iterative subgraph evaluation, question
782 generation, and answer validation. In addition, our benchmark includes baseline model
783 evaluations involving multiple LLMs and LLM-based KG-RAG models.
784

Appendices

785

786	A Related Works	19
787	A.1 KG-RAG Methods.	19
788	A.2 KGQA Benchmarks.	20
789	B KGQA Datasets and Inspection Protocol	21
790	B.1 Basic Dataset Statistic	21
791	B.2 Recent KG-RAG Publications with WebQSP and/or CWQ	22
792	B.3 Manual Inspection Protocol	22
793	C Pitfalls of Existing KGQA Benchmarks	24
794	C.1 Current KGQA Dataset Issues	24
795	C.2 Limitations of Exact-Match Evaluation	27
796	D Prompt Templates and Generation Examples	28
797	D.1 Question Generation Prompt	28
798	D.2 SPARQL Validation Prompt	30
799	E Ablation Study on KGQAGen Design	31
800	F KGQAGen-10k Analysis	32
801	F.1 KGQAGen-10k Statistics	32
802	F.2 A Case Study of KGQAGen-10k	32
803	G Experimental Details	35
804	G.1 Model Specifications	35
805	G.2 Beyond Exact Match: Introducing LASM	36
806	G.3 Experimental Setup	37

807 **A Related Works**

808 This section provides comprehensive technical details and broader context for the related work
809 summarized in Section 2, focusing on the methodological foundations and evaluation challenges that
810 motivate our framework.

811 **A.1 KG-RAG Methods.**

812 Knowledge graph-based retrieval-augmented generation systems address hallucination and factual
813 grounding limitations in large language models by integrating structured symbolic knowledge into
814 the generation process. These systems retrieve relevant subgraphs based on input queries to serve
815 as structured context, effectively combining broad parametric knowledge with the precision and
816 verifiability of knowledge bases.

817 Recent developments have produced several distinct architectural approaches addressing different
818 aspects of the integration challenge. Graph-guided reasoning methods, exemplified by Reasoning-on-
819 Graph (RoG) [33], enable interpretable multi-step reasoning by having language models explicitly

820 verbalize their traversal through knowledge graph structures, generating relation paths that are then
821 grounded in actual graph connections. This approach provides both answer generation and explanation
822 capabilities, making the reasoning process more transparent and debuggable. Extensions include
823 GNN-RAG [36], which incorporates graph neural networks to better capture structural patterns
824 in retrieved subgraphs. In contrast, constrained generation approaches like Graph-Constrained
825 Reasoning (GCR) [34] focus on ensuring faithful outputs by incorporating explicit graph-based
826 constraints into the decoding process, using KG-Trie structures to restrict generation to only those
827 paths that exist in the knowledge base.

828 Modular and memory-augmented architectures represent another important direction. Frameworks
829 like FRAG [18] propose adaptive combinations of different retrieval and generation components,
830 while Generate-on-Graph [72] treats the language model as both an agent and a knowledge graph
831 component, enabling interactive knowledge graph expansion during question answering. Memory-
832 augmented approaches such as MemQ [71] introduce dedicated memory modules that separate
833 language model reasoning from knowledge graph tool usage. More sophisticated integration strategies
834 include DeCAF [74], which jointly generates natural language answers and corresponding logical
835 forms, and ReknoS [66], which introduces abstract relationship reasoning through super-relations for
836 complex compositional queries.

837 **A.2 KGQA Benchmarks.**

838 Despite major progress in KGQA dataset construction, existing benchmarks exhibit persistent limi-
839 tations that hinder effective evaluation of modern KG-augmented retrieval systems. Early human-
840 curated datasets such as WebQuestions [6] and ComplexQuestions [5] focused on capturing authentic
841 user queries and introducing compositional constraints, but these resources are dominated by simple
842 factoid questions or contain ambiguities and incomplete annotations, making them insufficient for
843 testing models that require deeper reasoning and precise answer grounding.

844 To address these shortcomings, semi-automated and automated approaches have become prevalent.
845 Hybrid pipelines like LC-QuAD [62, 14] and GraphQuestions [54] use SPARQL templates and
846 graph-structured logic to guide question generation, with subsequent human editing for naturalness.
847 While this increases structural diversity and complexity, the template-based nature often leads to
848 unnatural or overly constrained question styles. More recent fully automated frameworks, including
849 Maestro [43], CHATTY-Gen [40], and DYNAMIC-KGQA [10], attempt to scale question generation
850 via rules, popularity heuristics, or dialogue simulation, yet these methods still struggle to balance
851 natural language fluency, logical completeness, and answer correctness for challenging multi-hop or
852 compositional queries.

853 The evolution of KGQA benchmarks reflects the field’s growing complexity, progressing from simple
854 factoid questions to sophisticated multi-hop reasoning scenarios. Early benchmarks established
855 foundational paradigms: WebQuestions pioneered collecting natural language questions through
856 search engine suggestions with human annotation, WebQuestionsSP [73] added SPARQL annotations
857 aligned with Freebase, and ComplexWebQuestions [58] advanced complexity by programmatically
858 generating SPARQL queries with logical constructs. However, these early datasets relied heavily
859 on Freebase, which ceased maintenance in 2016, motivating migration efforts to more sustainable
860 knowledge bases like DBpedia and Wikidata, though such migrations often introduced conceptual
861 mismatches and alignment difficulties. Specialized evaluation objectives drove targeted benchmark
862 development. GrailQA [20] introduced systematic evaluation of generalization scenarios including
863 i.i.d., compositional, and zero-shot settings, while KQAPro [8] emphasized compositional reasoning
864 through explicit program annotations and MetaQA [79] focused on multi-hop reasoning using con-
865 trolled movie domain knowledge graphs. Conversational benchmarks like CSQA [51] introduced
866 multi-turn dialogue structures, and recent work incorporated dialogue-style questions with corefer-
867 ence and ellipsis to better reflect real-world query patterns. Meanwhile, evaluation methodology
868 improvements from CBench [44] and SmartBench [42] provided comprehensive analysis of SPARQL
869 query complexity and linguistic diversity.

870 Despite these varied approaches, our systematic audit reveals that annotation quality remains a persis-
871 tent challenge across benchmarks, with even widely-used datasets like WebQSP and CWQ exhibiting
872 correctness rates below 60 percent. Furthermore, the predominant reliance on exact-match evalua-
873 tion metrics fails to capture semantic equivalence, leading to systematic underestimation of model

874 capabilities and underscoring the need for more rigorous benchmark construction methodologies that
875 prioritize both annotation accuracy and semantically-aware evaluation protocols.

876 **B KGQA Datasets and Inspection Protocol**

877 **B.1 Basic Dataset Statistic**

878 This section provides additional context for the dataset analysis presented in Section 3, reviewing
879 existing KGQA benchmarks and their construction methodologies. We categorize these datasets
880 based on their generation approaches and underlying knowledge bases. These datasets differ in
881 construction methods (manual vs. automated), target knowledge graphs (Freebase [7], DBpedia [3],
882 Wikidata [63], YAGO [55]), and reasoning complexity. Some emphasize compositional or multi-hop
883 reasoning, while others focus on multilingual support, dialogue capabilities, or logical form mapping.
884 Our review examines these datasets to understand their construction approaches, key features, and
885 relevance for evaluating knowledge graph-enhanced question answering systems.

886 **GraphQuestions** [54] is constructed through a semi-automated pipeline that begins with the gener-
887 ation of structured queries over Freebase. These queries are designed to capture varied reasoning
888 functions such as counting, superlatives, and conjunctions, as well as different structural complexities
889 and answer cardinalities. Queries that are infrequent or low-quality are filtered using web statistics.
890 The remaining set is then verbalized into natural language using pre-defined templates, creating a
891 dataset well-suited for evaluating semantic parsing and logical compositionality.

892 **WebQuestionsSP** [73] builds upon the WebQuestions [6] dataset by providing SPARQL annotations
893 aligned with Freebase. It employs a modular annotation interface that guides workers to select topic
894 entities, predicates, and relevant constraints, ensuring consistent semantic representation. The dataset
895 offers executable queries, enabling precise evaluation of models' ability to map natural language
896 questions to structured representations.

897 **ComplexWebQuestions** [58] extends WebQuestionsSP by programmatically generating more com-
898 plex SPARQL queries that incorporate logical constructs such as conjunctions, comparatives, and
899 superlatives. These queries are automatically translated into machine-generated questions, which
900 are then paraphrased into natural language by crowdworkers. The resulting dataset challenges QA
901 models with deeper compositional reasoning and varied surface forms.

902 **QALD-9** [39] is a multilingual benchmark constructed through manual curation, where annotators
903 write natural language questions and align them with SPARQL queries over DBpedia. It emphasizes
904 diversity in question types and supports multiple languages. Its extension, QALD-9-Plus, translates
905 these questions into additional languages and maps them to Wikidata, enabling cross-lingual and
906 cross-knowledge-base evaluation.

907 **MetaQA** [79] frames question answering as a latent variable problem, where both the topic entity and
908 the reasoning path are unobserved. Built over a movie-related knowledge graph, the dataset includes
909 1-hop, 2-hop, and 3-hop questions, designed to evaluate multi-step reasoning. A variational neural
910 framework is used to learn both entity disambiguation and graph-based reasoning simultaneously.

911 **SimpleDBpediaQA** [4] remaps the widely used SimpleQuestions dataset from Freebase to DBpedia.
912 This conversion involves aligning entities and predicates via owl:sameAs links and rewriting SPARQL
913 queries to account for conceptual mismatches such as directionality, ambiguity, and redirections.
914 The resulting dataset enables QA over an actively maintained knowledge base while preserving the
915 simplicity of the original benchmark.

916 **CSQA** [51] introduces a large-scale, multi-turn dialogue dataset for complex question answer-
917 ing over Wikidata. It combines crowd-sourced and in-house annotations to generate over 200K
918 question-answer pairs, including clarification, comparison, and logical reasoning within dialogue
919 context. CSQA is designed to benchmark conversational agents that require memory and contextual
920 understanding.

921 **LC-QuAD 1.0 and 2.0** [62, 14] are created through a three-stage semi-automated pipeline. First,
922 SPARQL queries are instantiated using templates and selected entities. Second, these queries are
923 mapped to question templates. Finally, crowdworkers paraphrase them into fluent natural language.
924 The datasets support a wide range of question types, including boolean, temporal, compositional, and
925 multi-relation queries, and target DBpedia and Wikidata respectively.

926 **FreeBaseQA** [23] compiles trivia-style factoid question–answer pairs from public sources and aligns
927 them with Freebase triples using a two-way entity linking approach. Human annotators then validate
928 relevance and correctness. The dataset features over 54K entity-answer pairs covering 28K unique
929 natural questions, offering high linguistic diversity and a challenging alternative to SimpleQuestions.

930 **Compositional Freebase Questions** [25] is developed to assess compositional generalization. It
931 begins by generating logical forms and corresponding SPARQL queries using a unified grammar.
932 Natural questions are written to express these logical forms. The dataset is then split using a
933 divergence-based strategy to maximize the gap in compositional structure between training and test
934 sets, enabling rigorous evaluation of generalization.

935 **GrailQA** [20] constructs question–answer pairs through a structured pipeline involving logical form
936 generation over Freebase, expert-written canonical questions, and crowd-sourced paraphrases. It
937 includes three evaluation splits: i.i.d., compositional, and zero-shot, making it suitable for analyzing
938 generalization across different reasoning complexities and linguistic expressions.

939 **QALD-9 Plus** [48] enhances QALD-9 by adding more questions, refining SPARQL annotations, and
940 expanding multilingual support. The updated version improves coverage of complex queries and
941 diverse linguistic phenomena, making it more suitable for modern multilingual QA systems.

942 **KQA Pro** [8] is constructed by generating compositional reasoning programs (KoPL) and correspond-
943 ing SPARQL queries over a curated knowledge base. These are paraphrased into natural questions
944 via crowdsourcing. The dataset supports fine-grained evaluation across logical reasoning categories
945 such as multi-hop inference, comparison, boolean logic, and temporal constraints.

946 **Dynamic-KGQA** [10] departs from static benchmarks by generating question–answer pairs on-the-
947 fly. It samples compact, semantically coherent subgraphs from YAGO 4.5 [55] and uses LLMs to
948 produce multi-hop questions grounded in these subgraphs. This design minimizes data leakage and
949 supports controlled, reproducible QA benchmarking with adaptive complexity.

950 **B.2 Recent KG-RAG Publications with WebQSP and/or CWQ**

951 WebQSP [73] and CWQ [58] have emerged as the primary benchmarks for empirical evaluation
952 in LLM-based knowledge graph question answering. Table 3 presents a chronological summary
953 of recent LLM-based KGQA models, indicating for each whether WebQSP and CWQ were used
954 for evaluation and providing a brief description of the model’s main approach. Notably, nearly 30
955 KG-RAG models released between 2022 and 2025 have adopted one or both datasets as central
956 components of their experimental protocols, despite the quality issues identified in Section 3.

957 Early LLM–KG hybrids, including ReaRev [35] and DECAF [74], set a precedent by reporting
958 results on both datasets, but typically focused on Hit@1 as the main metric. As the field progressed,
959 it became clear that Hit@1 alone could obscure over-generation and other qualitative issues. This
960 recognition prompted a shift in the community: subsequent agent-style models such as ToG [57],
961 RoG [33], ChatKBQA [32], and KD-CoT [65] began to report full precision, recall, and F1 scores
962 on both benchmarks, enabling more nuanced and meaningful comparison. By 2024 and 2025, eval-
963 uation on WebQSP and CWQ had become an established standard for the field. State-of-the-art
964 systems—such as GCR [34], Effi-QA [13], GNN-RAG [36], PoG [9], CLEAR-KGQA [67], and
965 ReKnoS [66]—universally benchmarked on these datasets before extending to additional corpora
966 or domain-specific tasks. Even in research targeting specialised knowledge graphs, models like
967 RARoK [75] and Efficient-G-Retriever [53] included WebQSP or CWQ for calibration and compar-
968 ability. As captured in Table 3, the adoption trajectory of these datasets not only reflects their ubiquity
969 but also highlights their role in shaping rigorous and transparent evaluation practices for the next
970 generation of KGQA systems.

971 **B.3 Manual Inspection Protocol**

972 We evaluate a broad selection of prominent KGQA benchmarks, including , and others commonly
973 used in recent KGQA literature. Each dataset provides a set of natural language questions paired
974 with ground-truth answers, and, where available, supporting triples or subgraphs from the underlying
975 knowledge graph (e.g., Freebase, Wikidata, DBpedia, WikiMovies and some subsets).

Table 3: Chronological summary of recent LLM-based KGQA models, with dataset usage based on reported experiments (✓).

Author [Citation]	Model	WebQSP	CWQ	Year	Short Introduction
Mavromatis et al. [35]	ReaRev	✓	✓	2022	LLM + GNNs refine reasoning on incomplete graphs.
Yu et al. [74]	DECAF	✓	✓	2022	Joint answer/logical form decoding from free-text retrieval.
Sun et al. [57]	ToG	✓	✓	2023	LLM agent explores KGs via beam search for deep, interpretable reasoning.
Luo et al. [33]	RoG	✓	✓	2023	Relation-grounded KG paths guide LLM reasoning with explanations.
Luo et al. [32]	ChatKBQA	✓	✓	2023	LLM-generated logical forms, improved with KG retrieval.
Wang et al. [65]	KD-CoT	✓	✓	2023	External KG knowledge injected into CoT reasoning.
Liu et al. [30]	DualR	✓	✓	2024	GNN for structural reasoning, frozen LLM for semantic reasoning.
Luo et al. [34]	GCR	✓	✓	2024	KG-Trie constrains LLM decoding for logic-faithful KG reasoning.
Dong et al. [13]	Effi-QA	✓	✓	2024	Iterative LLM planning, KG exploration, and self-reflection for QA.
Mavromatis et al. [36]	GNN-RAG	✓	✓	2024	GNN-based subgraph reasoning with LLM in RAG pipeline.
Xu et al. [70]	READS	✓	✓	2024	LLM decomposes KGQA into retrieval, pruning, inference.
Li et al. [26]	DoG	✓	✓	2024	LLM generates “well-formed chains” via constrained decoding.
Fang et al. [17]	KARPA	✓	✓	2024	LLM pre-plans, matches KG paths, reasons in training-free manner.
Xu et al. [72]	GoG	✓	✓	2024	LLM agent selects, generates, reasons on incomplete KGs.
Zhan et al. [75]	RARoK	✓	✓	2024	RAG-augmented CoT for complex medical KGQA.
Li et al. [27]	SubgraphRAG	✓	✓	2024	MLP + triple-scoring for efficient subgraph extraction.
Fang et al. [16]	DARA	✓		2024	LLM decomposes and grounds formal KG queries.
Hu et al. [22]	GRAG	✓		2024	Text-to-graph, retrieves/prunes subgraphs for RAG.
Xiong et al. [69]	Interactive-KBQA	✓	✓	2024	LLM agent generates SPARQL via multi-turn KB interaction.
Dehghan et al. [12]	EWEK-QA	✓	✓	2024	Web retrieval + KG triple extraction for citation-based QA.
Chen et al. [9]	PoG	✓	✓	2024	Self-correcting LLM planner for decomposed KGQA.
Wen et al. [67]	CLEAR-KGQA	✓	✓	2025	Interactive clarification and Bayesian inference for ambiguity.
Tan et al. [59]	Path-Over-Graphs	✓	✓	2025	LLM agent explores/prunes multi-hop KG paths.
Wang et al. [66]	ReKnoS	✓	✓	2025	Aggregates “super-relations” for LLM forward/backward reasoning.
Xu et al. [71]	MemQ	✓	✓	2025	Memory module separates LLM reasoning from KG tool use.
Gao et al. [18]	FRAG	✓	✓	2025	Modular KG-RAG adapts retrieval to query complexity.
Shen et al. [52]	RwT	✓	✓	2025	LLM-guided MCTS refines KG reasoning chains.
Solanki et al. [53]	Efficient-G-Retriever	✓		2025	Attention-based subgraph retriever for LLM-aligned RAG.
Tang et al. [60]	GGI-MAB	✓	✓	2025	Multi-armed bandit adapts RAG retrieval for KGQA.
Zhang et al. [78]	TrustUGA	✓		2025	Unified Condition Graph, two-level LLM querying.

976 For our systematic audit, we followed a unified sampling and review protocol. For WebQSP and
977 CWQ—the most widely adopted benchmarks—we randomly sampled 100 and 300 test examples,
978 respectively, to ensure sufficient coverage of prevalent error types. For all other datasets, we
979 sampled 60 examples per set to enable a balanced cross-dataset comparison. Each sampled item was
980 independently reviewed by two annotators with KGQA expertise, resolving disagreements through
981 discussion.

982 During inspection, we assessed each example along three main dimensions: (1) factual correctness
983 of the annotated answer; (2) clarity and appropriateness of the question; and (3) faithfulness of the
984 supporting SPARQL, where available. We flagged instances with incorrect, incomplete, or ambiguous
985 annotations, as well as questions that were underspecified, trivial, or unanswerable.

986 C Pitfalls of Existing KGQA Benchmarks

987 Many benchmarks are anchored to deprecated resources like Freebase, and even migration efforts
988 to Wikidata [46] introduce further inconsistencies in entity mappings and answer verification. Most
989 critically, almost all existing datasets are evaluated using rigid exact match (EM) metrics, which fail to
990 recognize semantically equivalent answers phrased differently. In summary, current KGQA datasets
991 face two central challenges: (1) data quality issues, including inaccurate, incomplete, or artificial
992 annotations, and (2) narrow evaluation protocols that do not capture true semantic correctness. These
993 limitations underscore the need for new benchmarks—such as KGQAGen—that emphasize both
994 annotation quality and robust, semantically-aware evaluation.

995 C.1 Current KGQA Dataset Issues

996 This appendix presents a detailed case study of data quality issues involved the most widely used
997 KGQA benchmarks, drawing from our broader review of 16 major datasets. For each dataset,
998 we randomly sampled question-answer pairs and performed careful manual verification following
999 the evaluation criteria in Section B. By presenting both problematic examples and the reasoning
1000 behind their classification, this analysis provides concrete, case-based evidence that complements the
1001 aggregate statistics reported in Table 1 and discussed in Section 3.

1002 Our review identifies three principal categories of data quality problems, as defined in Section 3.1:
1003 **Inaccurate Ground Truth Answers, Low-Quality or Ambiguous Questions, and Limitations**
1004 **of Exact-Match Evaluation.** These recurring issues—rooted in annotation, question design, and
1005 evaluation protocols—can seriously undermine the reliability of KGQA benchmarks. A compre-
1006 hensive breakdown, along with additional annotated examples for each dataset, is provided in the
1007 supplementary materials at the shared repository⁵.

1008 C.1.1 Inaccurate Ground Truth Answers

1009 Following the categorization presented in Section 3, we provide additional examples for each type of
1010 answer annotation error identified in our analysis: A major challenge in KGQA benchmarks is the
1011 prevalence of ground truth annotation errors, including incorrect, incomplete, or outdated answers.
1012 These undermine evaluation reliability and can mislead model training.

1013 **Incorrect annotations.** Incorrect annotations occur when the labeled answer does not actually
1014 address the question or contains factual mistakes.

- 1015 • **ID: WebQTest-273**
1016 **Question:** When did Michael Jordan return to the NBA?
1017 **Answer:** 1984
1018 **Issue:** 1984 is the year of his NBA debut, not his return. The correct answer should be 1995.
- 1019 • **ID: QALD9Plus-120**
1020 **Question:** Who is the daughter of Bill Clinton married to?
1021 **Answer:** Chelsea Clinton
1022 **Issue:** Answer is Chelsea Clinton herself, not her spouse. The correct answer should be
1023 Marc Mezvinsky.
- 1024 • **ID: GrailQA-2101221015000**
1025 **Question:** The 1912–13 Scottish Cup season is part of what sports league championship?
1026 **Answer:** 1913 Scottish Cup Final
1027 **Issue:** The answer refers to a final match, not the correct league entity (“Scottish Cup”).
- 1028 • **ID: DynamicKGQA-26720**
1029 **Question:** Which actor starred in both ‘The Hospital’ and was born in the same country as
1030 Albert Einstein?
1031 **Answer:** George C. Scott
1032 **Issue:** George C. Scott was born in the United States, while Einstein was born in Germany.
1033 The answer does not satisfy the nationality constraint.

⁵https://drive.google.com/drive/folders/1hH-NxqbUk0SeLC3q1ifLpq6i01Byait4?usp=drive_link

1034 **Outdated answers.** Outdated answers reflect facts that were once correct but have become obsolete
1035 due to real-world changes and key issue is the outdated knowledge source.

- 1036 • **ID: WebQTest-182**
1037 **Question:** Who is Khloe Kardashian’s husband?
1038 **Answer:** Lamar Odom
1039 **Issue:** The answer is outdated; Khloe Kardashian and Lamar Odom divorced in 2016.
- 1040 • **ID: CWQ-1182_ac67410188d0f2258139a3c84773885e**
1041 **Question:** Who is the current queen of the location whose time zone is Central Western
1042 Time Zone?
1043 **Answer:** Elizabeth II
1044 **Issue:** The answer is outdated. The Central Western Time Zone (UTC+8:45) corresponds
1045 to a small region in Australia. Since Queen Elizabeth II passed away in 2022, the current
1046 monarch of Australia is King Charles III.
- 1047 • **ID: KQAPro-14**
1048 **Question:** What city’s population is 6,690,432?
1049 **Answer:** Dalian
1050 **Issue:** The question lacks a specified time period, making the answer potentially outdated.
1051 Additionally, there is no country constraint, which may lead to ambiguity.
- 1052 • **ID: FreebaseQA-eval-836**
1053 **Question:** Who is currently the creative director at the house of Chanel?
1054 **Answer:** Karl Lagerfeld
1055 **Issue:** Karl Lagerfeld passed away in 2019 and no longer holds this position.
- 1056 • **ID: ComplexQuestions-33**
1057 **Question:** Who is Greece’s leader now?
1058 **Answer:** Karolos Papoulias
1059 **Issue:** The answer is outdated; Karolos Papoulias served as president until 2015. The correct
1060 answer should be the current leader.
- 1061 • **ID: MetaQA-46**
1062 **Question:** What were the release dates of films directed by the director of [Rosemary’s
1063 Baby]?
1064 **Answer:** 1986, 1948, 1992, 1982, 1994, 1979, 1999, 1965, 1974, 1967, 1971, 2002, 1988,
1065 2005, 1976, 2011, 2010, 2013
1066 **Issue:** The answer is incomplete and outdated; it does not include the director’s most recent
1067 films, such as a 2023 release.

1068 **Incomplete annotations.** Incomplete annotations occur when the gold set omits other valid answers,
1069 penalizing models that provide equally correct alternatives.

- 1070 • **ID: CWQ-124_7360e892294860c6ef7ad9a10e540e1b**
1071 **Question:** What movie did the author who published editions for *Notes from My Travels*
1072 direct?
1073 **Answer:** Unbroken
1074 **Issue:** The annotation is incomplete. The book *Notes from My Travels* was written by
1075 Angelina Jolie, who has directed multiple films, including *In the Land of Blood and Honey*
1076 (2011), *Unbroken* (2014), *By the Sea* (2015), and *First They Killed My Father* (2017).
- 1077 • **ID: GrailQA-2100689004000**
1078 **Question:** What fossil specimen dates from the Eocene?
1079 **Answer:** *Darwinius masillae*
1080 **Issue:** The annotation is incomplete; many Eocene fossils (e.g., *Moeritherium lyonsi*) would
1081 also be valid answers.
- 1082 • **ID: DynamicKGQA-14927**
1083 **Question:** Which American professor at the University of Wisconsin–Madison worked in
1084 the same city where Charles V. Bardeen died?
1085 **Answer:** E. Ray Stevens
1086 **Issue:** The annotation is incomplete; any American professor at UW–Madison would satisfy
1087 the condition, not just E. Ray Stevens.

- 1088 • **ID: DynamicKGQA-24330**
- 1089 **Question:** Which athlete from Novosibirsk shares their nationality with the owner of the
- 1090 United Shipbuilding Corporation?
- 1091 **Answer:** Yevgeni Nikolayevich Andreyev
- 1092 **Issue:** The annotation is incomplete; multiple Russian athletes from Novosibirsk could be
- 1093 correct.
- 1094 • **ID: GraphQuestions-281000202**
- 1095 **Question:** The British coat of arms has which heraldic supporters included?
- 1096 **Answer:** Lion
- 1097 **Issue:** The annotation is incomplete; both a lion (left) and a unicorn (right) are supporters.
- 1098 Only listing “lion” is insufficient.
- 1099 • **ID: KQAPro-72547**
- 1100 **Question:** What person has a notable work titled "The Scorpion King," which was produced
- 1101 by Vince McMahon?
- 1102 **Answer:** Dwayne Johnson
- 1103 **Issue:** Incomplete annotation; any individual involved in the production could be considered
- 1104 correct, not just Dwayne Johnson.

1105 C.1.2 Low-Quality or Ambiguous Questions

1106 Low-quality questions include those that are ambiguous, overly simple, or fundamentally unanswer-
 1107 able. We highlight key cases by dataset.

1108 **Ambiguous phrasing.** Ambiguous questions make it unclear what the intended answer should be,
 1109 undermining evaluation objectivity.

- 1110 • **ID: GrailQAPlus-2101990009000**
- 1111 **Question:** Which automotive designer designed NA?
- 1112 **Answer:** Koichi Hayashi, Tom Matano, Bob Hall
- 1113 **Issue:** The meaning of “NA” is ambiguous and could refer to multiple models or entities.
- 1114 • **ID: GraphQuestions-358000401**
- 1115 **Question:** sci came after what video game engine?
- 1116 **Answer:** Adventure Game Interpreter
- 1117 **Issue:** The question is ambiguous as “sci” does not have a defined meaning or Wikidata
- 1118 label.
- 1119 • **ID: GrailQA-3200295001000**
- 1120 **Question:** What is the number of theater plays that are in mystery?
- 1121 **Answer:** 1
- 1122 **Issue:** The question is too vague and does not provide enough context to determine the
- 1123 correct answer.
- 1124 • **ID: WebQTest-108**
- 1125 **Question:** What was the book written by Charles Darwin?
- 1126 **Answer (length: 153, unique: 94):** On evolution, The Autobiography of Charles Darwin,
- 1127 The Voyage of the Beagle, The Origin of Species, ...
- 1128 **Issue:** The question is ambiguous—Darwin wrote many books, but the prompt refers to
- 1129 “the” book. The ground truth annotation is also excessively broad and noisy, containing 153
- 1130 answers (including duplicates), which undermines reliable evaluation.

1131 **Low-complexity questions.** Low-complexity questions can be solved by simple lookup or direct
 1132 retrieval, offering little test of reasoning or multi-hop capabilities.

- 1133 • **ID: QALD9Plus-143**
- 1134 **Question:** What is the area code of Berlin?
- 1135 **Answer:** 030
- 1136 **Issue:** This is a straightforward 1-hop fact lookup with no reasoning required.
- 1137 • **ID: FreebaseQA-eval-1536**
- 1138 **Question:** Who wrote the 1970 book ‘Future Shock’?

1139 **Answer:** Alvin Toffler
1140 **Issue:** This is a basic factoid query that tests only surface-level knowledge.

1141 • **ID: FreebaseQA-eval-2363**
1142 **Question:** In which ocean is the island of Madeira?
1143 **Answer:** Atlantic Ocean
1144 **Issue:** This is a simple fact lookup with no reasoning challenge.

1145 • **ID: CSQA-7**
1146 **Question:** Which sex does Wolfgang Brandstetter belong to?
1147 **Answer:** male
1148 **Issue:** This is a 1-hop lookup question with no requirement for reasoning ability.

1149 • **ID: SimpleDBpedia-17597**
1150 **Question:** What was Karl Dönitz’s place of death?
1151 **Answer:** Schleswig Holstein
1152 **Issue:** This is a simple factual retrieval, requiring no reasoning.

1153 **Unanswerable, subjective, or ill-formed questions.** Such questions may rest on false premises,
1154 omit crucial context, or rely on subjective interpretations.

1155 • **ID: GrailQA-2104467004000**
1156 **Question:** The time zone from UTC of 12.75 has been offset what number of times?
1157 **Answer:** 1
1158 **Issue:** The question is uninterpretable; UTC+12.75 is not a standard offset, and the phrasing
1159 lacks clarity.

1160 • **ID: QALD9Plus-81**
1161 **Question:** Butch Otter is the governor of which U.S. state?
1162 **Answer:** [Missing Ground Truth]
1163 **Issue:** Unanswerable in the present; Butch Otter no longer holds that position.

1164 • **ID: FreebaseQA-eval-3290**
1165 **Question:** What is measured in Scoville units?
1166 **Answer:** Pungency
1167 **Issue:** Subjective; the question could accept “spiciness” or “pungency,” but only one is
1168 annotated as correct.

1169 • **ID: GraphQuestions-462000201**
1170 **Question:** Find the bearers of the coat of arms granted by queen.
1171 **Answer:** Western Australia
1172 **Issue:** The question does not specify which queen or which coat of arms, making it
1173 ambiguous and unanswerable.

1174 • **ID: KQAPro-3**
1175 **Question:** Which has lower elevation above sea level, Bristol or Jerusalem whose ISNI is
1176 0000 0001 2158 6491?
1177 **Answer:** Bristol
1178 **Issue:** Problematic: the Jerusalem referenced is a musician, not a location. Multiple cities
1179 named Bristol exist, with no way to determine which is intended.

1180 C.2 Limitations of Exact-Match Evaluation

1181 Existing KGQA benchmarks are further limited by their reliance on rigid exact-match evaluation
1182 protocols. Such criteria do not accommodate semantically correct answers that are phrased differently
1183 from the annotated ground truth. As a result, models are often penalized for generating correct
1184 answers that differ only in surface form, leading to false negatives and an underestimation of true
1185 model performance.

1186 • **ID: QALD9-174**
1187 **Question:** Who is the novelist of the work a song of ice and fire?
1188 **Ground Truth:** George_R._R._Martin
1189 **Issue:** Other semantically correct forms such as "George Raymond Richard Martin,"
1190 "George R. R. Martin" (with or without punctuation), "G. R. R. Martin," "George RR

1191 Martin," "Martin, George R. R.," "George R. Martin," or "G.R.R. Martin" are equally valid.
 1192 Exact-match evaluation penalizes correct answers that differ in surface form or formatting.

- 1193 • **ID: QALD9-114**
 1194 **Question:** How big is the earth's diameter?
 1195 **Ground Truth:** 1.2742e+07
 1196 **Issue:** Acceptable answers include "12,742 km," "12,742,000 meters," "about 7,918 miles,"
 1197 "1.2742 × 10⁷ meters," and "approximately 12,700 kilometers." Variations in units, notation,
 1198 or approximation are all reasonable, but exact match evaluation may reject them as incorrect.
- 1199 • **ID: CSQA-60**
 1200 **Question:** Which nucleic acid sequence encodes Ufm1-specific protease 2?
 1201 **Ground Truth:** Ufsp2
 1202 **Issue:** Other valid forms include "Ufsp2 gene," "the gene encoding Ufm1-specific protease
 1203 2," "gene symbol: UFSP2," or Ensembl/NCBI identifiers. Exact match allows only the
 1204 annotated form, penalizing equally correct alternatives.
- 1205 • **ID: WebQTest-6**
 1206 **Question:** Where is JaMarcus Russell from?
 1207 **Ground Truth:** Mobile
 1208 **Issue:** Answers such as "Mobile, Alabama," "the city of Mobile," "Mobile (city)," "Mobile,
 1209 AL," or "JaMarcus Russell was born in Mobile, Alabama" all convey the same information,
 1210 but may not be accepted unless they match the ground truth exactly.
- 1211 • **ID: FreebaseQA-eval-607**
 1212 **Question:** Who wrote the 1990 Booker Prize winner Possession?
 1213 **Ground Truth:** a. s. byatt
 1214 **Issue:** Other correct answers like "Antonia Susan Byatt," "Dame A. S. Byatt," or "Antonia
 1215 Byatt" are semantically equivalent, but only the annotated form is accepted under exact
 1216 match.

1217 D Prompt Templates and Generation Examples

1218 This section provides detailed documentation of the prompt templates used in KGQAGen, along with
 1219 representative examples of generated questions and error cases. We first present the core generator
 1220 prompt that guides question construction and knowledge sufficiency checking. We then describe the
 1221 simplified evaluator prompt used during answer validation. Finally, we analyze example outputs and
 1222 common error patterns to illustrate the framework's capabilities and limitations.

1223 D.1 Question Generation Prompt

1224 The generator component of KGQAGen utilizes a carefully designed prompt template to ensure high-
 1225 quality question generation. The prompt consists of input specifications and structured generation
 1226 rules that guide the LLM in producing well-formed question-answer instances.

1227 The input format specifies RDF triples from Wikidata, where each triple contains an entity label with
 1228 its Q-ID, a predicate label with its P-ID, and another entity label with Q-ID. The LLM evaluates
 1229 whether the given subgraph provides sufficient information for generating a meaningful KGQA
 1230 instance. The generation rules enforce five key requirements for question construction. (1) Reasoning
 1231 complexity demands that generated questions involve at least 2-hop logical reasoning paths. The
 1232 framework prioritizes specific, instance-level entities while avoiding generic categories. Factual
 1233 constraints are incorporated only when necessary for disambiguation or meaningful answer space
 1234 reduction. (2) Entity selection criteria require that candidate entities must be concrete instances
 1235 rather than abstract types. The system favors entities that enable meaningful multi-hop paths through
 1236 affiliations, awards, locations, or temporal relationships, while avoiding generic class-level concepts.
 1237 (3) Question difficulty ensures that questions are designed to require structured knowledge graph
 1238 reasoning by incorporating inverse relations, numerical constraints, comparative logic, or set-based
 1239 conditions. The framework employs factual filters to reduce ambiguity and ensures questions cannot
 1240 be answered through general knowledge alone. (4) Natural language quality mandates that questions
 1241 must use natural, fluent phrasing typical of real user queries. The prompt enforces self-contained and
 1242 concise formulation while avoiding references to underlying data structures or unnecessary repetition.
 1243 (5) Semantic clarity requires that all generated questions must unambiguously specify their intended

1244 answers. The prompt explicitly prohibits vague or underspecified formulations that could lead to
1245 multiple valid interpretations.

1246 The LLM returns a JSON response indicating either insufficiency with candidate entities for expansion
1247 or a complete question-answer instance containing the natural language question, answer set,
1248 supporting proof triples, and corresponding logical constraints. The full prompt template is provided
1249 below:

Prompt

```
You are given a small set of RDF triples from Wikidata.
Format: Each triple is a 3-item array: ["<label> (<Q-ID>)", "<predicate> (<P-ID>)", "<label> (<Q-ID>)"]
Triples: {triples}
Your task is to determine whether this subgraph is sufficient to support a challenging and non-trivial question for a knowledge graph question answering (KGQA) benchmark.
Guidelines:
1. Reasoning Depth
  • Prefer questions requiring at least 2-hop reasoning.
  • Avoid generic topics or subclass chains-focus on instance-level, specific entities.
  • Use factual constraints (e.g., date, affiliation) only when needed to disambiguate the answer or add meaningful specificity.
  • Do not over-constrain-include only what is necessary to yield a specific answer.
2. Entity Selection and Expansion
  • Focus on concrete, instance-level entities (e.g., Q7186), not types like Q5 (human) or Q11424 (film).
  • Avoid generic classes like "scientist", "award", or "event", and relations like "subclass of" or "instance of".
  • Prefer entities and paths supporting deeper reasoning-e.g., affiliations, recognitions, or spatiotemporal links.
3. Difficulty
  • Encourage inverse relations, comparative logic, date/number filters, or set membership.
  • Ensure the answer cannot be derived from general knowledge alone.
  • The subgraph must contain all supporting information to answer the question.
4. Naturalness
  • Phrase the question as a fluent, self-contained query a user might ask.
  • Avoid references to the input format (e.g., "triples", "given data").
  • Do not use phrases like:
    - "from the given data"
    - "among these entities"
    - "listed here"
5. Clarity
  • The question must be unambiguous and logically imply a unique, specific answer.
  • Avoid vague or underspecified language.

Output Format:
If the graph is not sufficient, return: { "sufficient": false, "candidate": [<QID>, ..., <QID>] }
If sufficient, return: { "sufficient": true, "question": "<natural-language question>", "answer": ["<answer-label (QID)>", "..."], "proof": [ ["<label (QID)>", "<predicate (PID)>", "<label (QID)>"], ... ] }
Return strict JSON only - no commentary.
```

1250

Few-shot Example

```
Example 1: Triples: [ ["Johann Martin Schleyer (Q12712)", "nominated for (P1411)", "Nobel Peace Prize (Q35637)", ["International Volapöfck Academy (Q3358168)", "founded by (P112)", "Johann Martin Schleyer (Q12712)"], ["Johann Martin Schleyer (Q12712)", "place of birth (P19)", "Oberlauda (Q885402)"] ] ]
Output: { "sufficient": true, "question": "Who among the nominees for the Nobel Peace Prize was also the founder of International Volapöfck Academy?", "answer": ["Johann Martin Schleyer (Q12712)", "proof": [ ["Johann Martin Schleyer (Q12712)", "nominated for (P1411)", "Nobel Peace Prize (Q35637)", ["International Volapöfck Academy (Q3358168)", "founded by (P112)", "Johann Martin Schleyer (Q12712)"] ] ] }
```

```
Example 2: Triples: [ ["Karakalpakstan (Q484245)", "capital (P36)", "Nukus (Q489898)", ["Karakalpakstan (Q484245)", "shares border with (P47)", "Mangystau Region (Q238931)", ["Karakalpakstan (Q484245)", "official language (P37)", "Karakalpak (Q33541)", ["Karakalpakstan (Q484245)", "country (P17)", "Uzbekistan (Q265)"] ] ] ]
Output: { "sufficient": false, "candidate": ["Q33541", "Q489898", "Q238931"] }
```

```
Example 3: Triples: [ ["Astronomy and Astrophysics (Q752075)", "publisher (P123)", "EDP Sciences (Q114404)", ["Astronomy and Astrophysics (Q752075)", "editor (P98)", "Thierry Forveille (Q46260676)", ["Zeitschrift für Astrophysik (Q3575110)", "followed by (P156)", "Astronomy and Astrophysics (Q752075)"] ] ] ]
Output: { "sufficient": true, "question": "What astronomical journal, published by EDP Sciences and edited by Thierry Forveille, succeeded Zeitschrift für Astrophysik as its immediate follower?", "answer": ["Astronomy and Astrophysics (Q752075)", "proof": [ ["Astronomy and Astrophysics (Q752075)", "publisher (P123)", "EDP Sciences (Q114404)", ["Astronomy and Astrophysics (Q752075)", "editor (P98)", "Thierry Forveille (Q46260676)", ["Zeitschrift für Astrophysik (Q3575110)", "followed by (P156)", "Astronomy and Astrophysics (Q752075)"] ] ] ] }
```

1251

1252 This structured prompt design ensures consistent generation of high-quality, diverse, and well-
1253 specified question-answer pairs for the benchmark dataset. The prompt template balances multiple
1254 objectives: maintaining reasoning complexity, ensuring natural language quality, and guaranteeing
1255 answer specificity. By enforcing these requirements through explicit rules and format specifications,
1256 we enable systematic generation of challenging yet well-formed KGQA instances.

1257 D.2 SPARQL Validation Prompt

1258 The validation component uses a focused prompt template for verifying and refining SPARQL queries.
1259 This streamlined prompt specifically addresses query correction when execution results differ from
1260 intended answers, ensuring both syntactic correctness and semantic alignment.

1261 The validation process operates on the principle of iterative refinement. When a generated SPARQL
1262 query fails to return expected results or returns empty result sets, the validation component engages a
1263 lightweight language model to diagnose and correct the query. This approach recognizes that initial
1264 query generation may suffer from syntactic errors, incorrect entity identifiers, or inappropriate query
1265 structure that prevents successful execution against the Wikidata endpoint.

1266 The prompt template emphasizes simplicity and executability in query revision. Rather than attempt-
1267 ing complex query transformations, the validation component focuses on ensuring that revised queries
1268 use only essential triple patterns and avoid unnecessary complexity that might introduce additional
1269 failure points. The template specifically discourages the use of optional clauses and filter conditions
1270 unless they are strictly necessary for answering the question, as these constructs often lead to query
1271 execution failures or unexpected empty results. Additionally, the validation process enforces struc-
1272 tural requirements that ensure compatibility with the Wikidata query service. All revised queries must
1273 terminate with a single SERVICE wikibase:label clause to retrieve human-readable English labels
1274 for entities, maintaining consistency with the expected output format. The prompt also mandates
1275 syntactic validity and direct executability at the official Wikidata SPARQL endpoint, ensuring that

1276 corrected queries can be verified immediately. The output format maintains strict JSON formatting
 1277 requirements to facilitate automated processing. The validation component returns only the corrected
 1278 SPARQL query without additional commentary or explanation, enabling seamless integration into
 1279 the broader validation pipeline. This focused approach allows for rapid iteration and correction when
 1280 initial query generation produces non-executable or semantically misaligned queries.

1281 Through this validation mechanism, KGQAGen ensures that all retained question-answer pairs are
 1282 grounded in verifiable SPARQL queries that can be executed against the knowledge base. This
 1283 constraint provides a strong guarantee of answer correctness and enables ongoing validation as the
 1284 underlying knowledge graph evolves over time.

```

Prompt
You are given a SPARQL query over Wikidata that returned no results.
Question: {question}
Original SPARQL: {sparql}
Your task is to revise the query so that it returns valid results from Wikidata.
Revision Guidelines:
• Use only essential triple patterns. Avoid OPTIONAL and FILTER clauses unless
  strictly necessary.
• The query must end with a single SERVICE wikibase:label clause to retrieve
  English labels.
• Ensure the query is syntactically valid and directly executable at https:
  //query.wikidata.org.
Output Format: Return a single JSON object in the exact format below - no
commentary, no markdown:
{
  "correct_sparql": "<REVISED SPARQL QUERY HERE>"
}
  
```

1285
 1286 This structured prompt design ensures consistent generation of high-quality, diverse, and well-
 1287 specified question-answer pairs for the benchmark dataset. The prompt template balances multiple
 1288 objectives by maintaining reasoning complexity, ensuring natural language quality, and guaranteeing
 1289 answer specificity. By enforcing these requirements through explicit rules and format specifications,
 1290 we enable systematic generation of challenging yet well-formed KGQA instances.

1291 E Ablation Study on KGQAGen Design

1292 **Setup.** We perform an ablation study to quantify the contribution of core components in the
 1293 KGQAGen-10k generation pipeline. In our framework, questions are generated through iterative
 1294 LLM-guided subgraph expansion combined with symbolic SPARQL-based verification. To isolate
 1295 the effect of each component, we consider three configurations: (A1) random subgraph selection
 1296 replacing LLM-guided expansion, (A2) alternative generation LLMs, and (A3) inclusion or exclusion
 1297 of SPARQL-based answer validation. A3 is treated as a cross-setting axis applied to all variants. We
 1298 randomly sample 100 seed entities and report two metrics: *generation success rate* (percentage of
 1299 seeds producing valid questions) and *end-to-end success rate* (percentage producing fully validated
 1300 QA pairs).

Setting	Generation Rate (%)	End-to-End Rate (%)	Validation Effect (%)
Full KGQAGen (GPT-4.1)	94	77	+17
A1: Random Subgraph	100	62	+15
A2-1: LLaMA-3-70B	57	35	+4
A2-2: GPT-4o-mini	65	41	+6

Table 4: Ablation results on 100 randomly selected seed entities. “Validation Effect” denotes the improvement in valid QA pairs after applying SPARQL verification.

1301 **Findings.** The full KGQAGen pipeline achieves a 94% generation success rate and a 77% end-to-end
 1302 validation rate. Removing LLM-guided subgraph selection (A1) maintains full generation coverage

1303 but lowers validation to 62%, as random sampling often introduces extraneous entities that cause
 1304 ambiguity. To assess the effect of model choice, we replace GPT-4.1 with smaller LLMs under the
 1305 same setup. LLaMA-3-70B generates 57 questions with 35 validated (61% of generated), while
 1306 GPT-4o-mini produces 65 questions with 41 validated (63% of generated). Smaller models frequently
 1307 simplify multi-hop reasoning to single-hop relations or omit relational constraints. In contrast, the
 1308 full KGQAGen pipeline with GPT-4.1 yields 82 generated and 77 validated questions, achieving
 1309 both higher coverage and correctness. These results demonstrate that model capacity and structure-
 1310 aware subgraph guidance jointly determine data quality, while SPARQL-based symbolic verification
 1311 mitigates hallucinations and constrains potential biases.

1312 To better understand the failure modes of symbolic verification, we analyzed 420 instances filtered
 1313 out by the SPARQL-based validation stage. The main error types include: empty query results (166
 1314 cases, 39.5%), excessively large answer sets (184 cases, 43.8%), and answer mismatches where the
 1315 annotated answer was absent from the SPARQL results (70 cases, 16.7%). These issues primarily
 1316 stem from incomplete or ambiguous question–answer pairs, limited KG coverage, or challenges
 1317 in SPARQL query formulation. Future work will strengthen the validation module with tighter
 1318 constraints, clarification prompts, and finer-grained error classification to improve explainability and
 1319 systematically address validation failures.

1320 **F KGQAGen-10k Analysis**

1321 **F.1 KGQAGen-10k Statistics**

1322 **Structural Complexity Metrics.** Following DyVal [81], we measure six structural aspects of the
 1323 supporting subgraphs for all 10,787 examples in KGQAGen-10k: (1) number of nodes, (2) number of
 1324 edges, (3) average degree, (4) reasoning depth (hops), (5) width (maximum frontier size), and (6)
 1325 extra links (non-chain connections). These metrics collectively capture both the scale and topological
 1326 difficulty of multi-hop reasoning paths.

Measure	Distribution Range	Percentage of Examples
Depth (hops)	2–5	98%
Nodes (entities)	5–30	84%
Edges (relations)	4–28	83%
Width (max frontier)	4–20	79%
Reachability	Fully connected	92%

Table 5: Graph-based reasoning complexity analysis for KGQAGen-10k following DyVal-style metrics. The dataset exhibits diverse and non-trivial structural properties, demonstrating its suitability for evaluating reasoning over complex knowledge graphs.

1327 Across these dimensions, most questions involve 2–5-hop reasoning chains over compact yet densely
 1328 connected subgraphs. The high connectivity (92%) and wide range of entity and relation counts
 1329 illustrate that KGQAGen-10k contains structurally rich and diverse reasoning contexts, providing a
 1330 challenging evaluation setting for KG-RAG and KGQA models.

1331 **F.2 A Case Study of KGQAGen-10k**

1332 An audit of 300 randomly sampled question–answer pairs from the entire 10,787-instance
 1333 KGQAGen-10k revealed 11 defective cases shown Table 6, a rate of error of 3.6%. Although this
 1334 figure is relatively low, these instances expose recurring weaknesses in the generation and verification
 1335 pipeline that warrant attention. The issues fall into three broad categories: self-answering prompts,
 1336 hallucinated or incomplete relations, and errors inherited from the source knowledge graph.

1337 The first issue, **self-answering questions**, was evident in items 4555 and 6931, where the question
 1338 text directly includes the target answer. For example, asking about a subclass that ‘has the same
 1339 meaning as the English-language word ‘city’ leaves little ambiguity about the expected answer.
 1340 Because our current verification process only checks that the SPARQL query returns at least one
 1341 result overlapping the proposed answer set, it overlooks this form of lexical leakage and accepts the
 1342 examples as valid.

1343 The second and more prevalent category involves **hallucinated or incomplete knowledge**. In six
1344 cases (IDs 8318, 1105, 1164, 1529, 1825, and 10469), the model generated questions based on
1345 nonexistent or incoherent relationships, such as attributing architectural roles to historical political
1346 figures. In these situations, the LLM still produces syntactically valid queries, sometimes by leverag-
1347 ing loosely related property paths, allowing the verifier’s overlap check to pass despite clear semantic
1348 errors. In a complementary failure mode, item 2297 exemplifies incomplete annotations: While
1349 multiple mathematicians satisfy the described criteria, only one is listed in the answer set. Since the
1350 verifier stops when it finds a match, it does not detect incompleteness.

1351 A third source of error originates not from the generation process but from the **underlying knowledge**
1352 **graph itself in Wikidata [63]**. Items 9572 and 2046 illustrate this point: one references an entity
1353 mistyped as a product, the other includes an unlabeled identifier. Our pipeline implicitly treats
1354 Wikidata typing and labeling as authoritative, so these issues remain undetected unless caught during
1355 manual review.

1356 The major cause of these verification failures lies in the limited scope of the current safeguard. The
1357 answer verification and refinement (detailed in Section 4.3) of our KGQAGen checks are whether the
1358 SPARQL query compiles and whether its result set overlaps the proposed answer. Although effective
1359 in filtering out broken or irrelevant queries, this approach does not account for key semantic and
1360 structural issues. It does not detect leakage of lexical answers, does not require precise predicate
1361 alignment, does not check the joint coherence of query constraints, and does not validate type or label
1362 accuracy within the knowledge graph.

1363 To address these gaps and further improve the quality of the data set beyond the current validation
1364 rate of 96.3%, we plan a series of targeted upgrades. First, we will implement a lexical filter to
1365 reject questions that contain their own answers. Second, we will enforce stricter predicate and type
1366 constraints within the generated queries to check against hallucinated relations. Third, we will
1367 introduce answer completeness audits through closure tests that ensure the full set of valid answers is
1368 captured. Finally, we will cross-check critical entity labels and types against alternative KG snapshots
1369 to catch inconsistencies and improve robustness across knowledge versions.

Table 6: KGQAGen-10k Error Analysis: Each case is displayed with concise issue diagnosis.

Field	Content
ID	4555
Question	What subclass of 'city or town' has both a GND ID and is identified as the same concept as the English-language word 'city'?
Answer	city
Issue	Self-answering: The question is trivial; the answer is explicitly stated in the question itself.
ID	8318
Question	Who is the architect of Estadio Nacional de Costa Rica that was also a successful candidate in multiple Costa Rican general elections?
Answer	Ricardo Jiménez Oreamuno
Issue	Hallucinated knowledge: The question implies an architect relationship that does not exist; Ricardo Jiménez Oreamuno was a politician, not an architect.
ID	1105
Question	Which person who has been an owner of the Shroud of Turin is not a human individual, but rather a dynastic house?
Answer	Geoffroi de Charny, House of Savoy, Jeanne de Vergy, pope
Issue	Hallucinated knowledge: Geoffroi de Charny and House of Savoy are individuals, not dynastic houses.
ID	1164
Question	Which person, who held citizenship in the Ming, Qing, and short-lived Zhou dynasties, was both the father of Wu Yingxiong and the spouse of both Chen Yuanyuan and Empress Zhang, and led a revolt known as the Revolt of the Three Feudatories?
Answer	Kangxi Emperor, Kingdom of Tungning, Qing dynasty, Wu Sangui, Zheng Jing
Issue	Hallucinated knowledge: Zheng Jing is not the spouse of Chen Yuanyuan.
ID	2297
Question	Which mathematician both lent his name to the principle maintained by WikiProject Mathematics and is explicitly named as its discoverer or inventor?
Answer	Johann Peter Gustav Lejeune Dirichlet
Issue	Hallucinated knowledge: Many mathematicians have principles named after them; the annotation is incomplete and not unique.
ID	1825
Question	Which concept, characterized as 'unusual', is named after 'unusual', is the main subject of an entity named after 'unusual', and is also cited as a partially coincident concept by 'rarity'?
Answer	frequency, scarcity, unusual
Issue	Hallucinated knowledge: None of the ground truths are the subject of "World's Weirdest Animals". Its main subject is "creature".
ID	10469
Question	Which musical artist is associated with the genre that is said to be the same as "vintage" and has also performed in the genre represented by 'retro style'?
Answer	Adam Tsarouchis, Ahmad Bersaudara, Anna Jantar, Gloomwood, Nina Shatskaya, Type-B, VCTR-SCTR, Wieczór na dworcu w Kansas City
Issue	Hallucinated knowledge: Wieczór na dworcu w Kansas City is not a musical artist, but a retro style song.
ID	1529
Question	Which artist created a work that depicts the astronomical event discovered by Pierre Gassendi, and has as its main subject the transit of Mercury?
Answer	Mercury Passing Before the Sun
Issue	Hallucinated knowledge: Mercury Passing Before the Sun is the artwork, not the artist. The artist is Giacomo Balla.
ID	6931
Question	Which type of underwater vehicle is classified as both a subclass of submersible and shares this property with bathyscaphe, narco-submarine, and Osprey-class submersible?
Answer	Osprey-class submersible, bathyscaphe, bathysphere, narco-submarine, submersible drilling rig
Issue	Self-answering: The answer is trivially the same as the question's subject; provides no substantive challenge.
ID	9572
Question	Which company has produced a product used for flow measurement that includes a flow meter as a part?
Answer	Sage Metering
Issue	Wikidata Mislabeled: The product of Sage Metering is labelled as "flow measurement" on Wikidata, but "flow measurement" is a task, not a product.
ID	2046
Question	Which tool that is a subclass of both 'physical tool' and is connected with 'level staff', is in turn a subclass of something that has the shape of a cylinder and is different from 'Rod'?
Answer	"Q9397141"
Issue	Wikidata Mislabeled: The entity "Q9397141" doesn't contain the natural language label.

1370 G Experimental Details

1371 This section provides comprehensive details about our experimental setup, including model configu-
1372 rations, training protocols, and evaluation procedures.

1373 G.1 Model Specifications

1374 We evaluate three categories of models on KGQAGen-10k: (1) **Pure Language Models**, (2) **KG-RAG**
1375 **Systems**, and (3) **LLMs with Supporting Graphs**. These categories reflect increasing levels of
1376 external knowledge integration, ranging from purely parametric reasoning to symbolic augmentation
1377 and perfect-evidence setups.

1378 **Pure Language Models** These models rely entirely on their internal parametric memory and are
1379 evaluated in a zero-shot setting, without any KG subgraph or retrieval. Their performance reflects
1380 inherent reasoning capability, factual coverage, and generalization.

- 1381 • **LLaMA-3.1-8B-Instruct** [19]: An 8B instruction-tuned model from Meta’s LLaMA 3.1
1382 series. It is optimized for following task-specific instructions and shows improved reasoning
1383 performance compared to earlier LLaMA versions.
- 1384 • **LLaMA2-7B** [61]: A general-purpose 7B model trained on publicly available data, serving
1385 as a foundational open-weight baseline for reasoning without instruction tuning.
- 1386 • **Mistral-7B-Instruct-v0.2** [2]: An instruction-following model based on Mistral-7B with
1387 a 32k context window and standard attention. It is designed for accurate and efficient
1388 long-context reasoning.
- 1389 • **GPT-4o-mini** [1]: A compact version of GPT-4o offering reduced latency and strong
1390 language understanding performance, suitable for real-time applications.
- 1391 • **GPT-4** [1]: OpenAI’s flagship model known for robust multi-step reasoning, long-context
1392 understanding, and generalization across a wide array of tasks.
- 1393 • **DeepSeek-Chat** [11]: A dialogue-oriented LLM developed by DeepSeek and fine-tuned for
1394 task completion and conversational fluency aligned with human feedback.
- 1395 • **GPT-4o** [1]: A unified multimodal model capable of handling text, image, and audio inputs.
1396 We use it in a text-only setup to assess its advanced reasoning capabilities.
- 1397 • **GPT-4.1** [41]: An updated variant of GPT-4 that improves long-context performance, factual
1398 grounding, and consistency in complex prompt execution.

1399 **KG-RAG Systems** Knowledge Graph Retrieval-Augmented Generation (KG-RAG) systems incor-
1400 porate structured symbolic evidence from a KG to assist reasoning and answer generation. These
1401 models access retrieved subgraphs at runtime and vary in how they integrate retrieved content—either
1402 as conditioning input or through decoding constraints.

- 1403 • **RoG (LLaMA2-7B)** [33]: Fine-tuned on KGQAGen-10k’s training split, using annotated sup-
1404 porting subgraphs to supervise faithful reasoning path generation for answer and explanation
1405 prediction.
- 1406 • **GCR (LLaMA-3.1 + GPT-4o)** [34]: Fine-tuned its path generation module with supporting
1407 subgraphs, leveraging a KG-Trie to constrain decoding and using GPT-4o for final answer
1408 synthesis.
- 1409 • **ToG (GPT-4o)** [57]: Adapted to Wikidata by replacing Freebase API calls with SPARQL
1410 queries and evaluated zero-shot, interactively exploring the KG and generating answers
1411 from the retrieved subgraph without fine-tuning or parameter adjustment.
- 1412 • **PoG (GPT-4o)** [9]: Applied as a zero-shot prompting-based agent that dynamically de-
1413 composes, explores, and self-corrects over Wikidata for each test question, without any
1414 dataset-specific fine-tuning.

1415 **LLM with Supporting Subgraph** To estimate the upper bound of KG-augmented QA performance,
 1416 we provide models with the gold supporting subgraph used during data generation. This simulates
 1417 a perfect-retrieval setting, where the model receives all and only the minimal evidence needed to
 1418 answer correctly. These experiments assess whether models can effectively reason over structured
 1419 KG input when retrieval is assumed to be ideal.

- 1420 • **LLaMA2-7B (w/ SP)** [61]: The model is provided with the gold subgraph and asked to
 1421 generate the answer. This tests the reasoning capacity of a smaller open-weight model under
 1422 ideal symbolic input.
- 1423 • **GPT-4o (w/ SP)** [1]: The same setup as above, but with GPT-4o as the base model. This con-
 1424 figuration reflects an upper-bound for KG-RAG systems when both retrieval and reasoning
 1425 are ideal.

1426 G.2 Beyond Exact Match: Introducing LASM

1427 While the limitations of exact match evaluation in KGQA are well-recognized [50, 64], few works
 1428 have proposed principled solutions. To address this gap, we introduce LLM-Assisted Semantic Match
 1429 (LASM), a novel evaluation scheme that goes beyond surface-level equivalence by leveraging the
 1430 semantic understanding capabilities of large language models.

1431 The core idea of LASM is to use an LLM verifier to assess semantic similarity between predicted
 1432 and ground truth answers. When a model’s prediction fails the exact string match, LASM invokes
 1433 a GPT-4o-mini judge to determine whether the prediction is semantically equivalent to the gold
 1434 answer. This approach enables LASM to properly credit models for generating meaningfully correct
 1435 responses that traditional metrics would overlook due to syntactic or lexical variation. To quantify the
 1436 impact of LASM, we compare model performance on the FreebaseQA dataset [23] under both exact
 1437 match and LASM evaluation. As shown in Table 7, LASM yields substantial improvements across all
 1438 key metrics, including accuracy (+5.3%), Hit@1 (+5.3%), and F1 (+5.0%). These gains demonstrate
 1439 the effectiveness of semantic matching in capturing valid model predictions that exact match misses.

Scoring	Accuracy	Hit@1	F1	Precision	Recall
Exact Match	90.39	90.39	88.08	87.08	90.39
LASM	95.72	95.67	93.12	92.04	95.65

Table 7: FreeBaseQA results with GPT-4o. LASM consistently recovers semantically correct predictions missed by exact match, leading to substantial metric improvements.

1440 Beyond offering a more robust and nuanced assessment of model outputs, LASM has important
 1441 implications for the development and evaluation of KGQA systems. By rewarding models for
 1442 semantic correctness rather than rigid string matching, LASM promotes the development of systems
 1443 that prioritize meaning over surface form. Moreover, as a fully automated method that does not
 1444 rely on dataset-specific rules or annotations, LASM is readily applicable to any KGQA benchmark,
 1445 enabling more meaningful cross-dataset comparisons.

1446 To ensure evaluation reliability, LASM is applied only when the prediction fails the exact match,
 1447 serving as a selective fallback rather than a replacement for EM. This design mitigates potential risks
 1448 of LLM-based scoring such as hallucinations or factual drift, as LASM is invoked for only a small
 1449 subset of cases. To assess its accuracy, we manually inspected 100 randomly sampled instances
 1450 where LASM was triggered. Two human reviewers independently verified semantic correctness
 1451 using Wikipedia, identifying only one false positive (1% error rate)—a case where the prediction
 1452 “*Austronesian languages*” was judged equivalent to the ground truth “*Oceanic*”. This low error rate
 1453 demonstrates that LASM provides a reliable supplement to exact match evaluation. Furthermore,
 1454 similar approaches in the NLG domain support the robustness of LLM-based evaluation, as shown by
 1455 recent studies such as G-Eval [31], FActScore [37], MT-Bench [80], and QAFactEval [68], which
 1456 report strong alignment between GPT-4-based judgments and human evaluation.

1457 In summary, LASM represents a principled and generalizable approach to overcoming the limitations
 1458 of traditional exact match evaluation in KGQA. By incorporating semantic awareness through LLM-
 1459 based similarity judgments, LASM provides a more reliable and nuanced assessment of model
 1460 performance, paving the way for the development of more robust question answering systems. As

1461 we will demonstrate through extensive experiments in Section 5, LASM offers a valuable tool for
1462 evaluating and advancing the state of the art in KGQA.

1463 **G.3 Experimental Setup**

1464 We evaluate all models on KGQAGen-10k using a standardized split of 8,629/ 1,079/1,079 train/e-
1465 val/test. For KG-RAG systems, we adapt each model to work with Wikidata by replacing their
1466 original knowledge base interfaces with SPARQL queries to the Wikidata endpoint.

1467 **Training and Inference Protocols** RoG [33] employs a planning-retrieval-reasoning pipeline
1468 where LLaMA2-7B first generates candidate relation paths, which are then matched against the
1469 knowledge graph using constrained breadth-first search. The retrieved reasoning paths, combined
1470 with the original question, guide the model to generate both answers and explanations. We fine-tune
1471 the entire pipeline on our training split, using the supporting subgraphs from dataset construction as
1472 supervision for faithful path generation.

1473 GCR [34] enforces graph faithfulness through constrained decoding. Prior to inference, we construct
1474 a KG-Trie index that efficiently captures all valid reasoning paths within a fixed hop limit. During
1475 generation, a fine-tuned LLaMA-3.1 model produces candidate paths under strict KG-Trie constraints,
1476 ensuring only valid graph traversals. These candidates are then passed to GPT-4o for inductive
1477 reasoning and answer synthesis. Similar to RoG, we leverage our supporting subgraphs for training
1478 the path generation component.

1479 In contrast, ToG [57] and PoG [9] operate without fine-tuning, treating the LLM as an agent that
1480 interactively explores the knowledge graph. ToG constructs reasoning trees by iteratively selecting
1481 relations and entities based on question semantics, while PoG enhances this with adaptive planning
1482 and self-correction mechanisms. For both models, we implement direct Wikidata integration, allowing
1483 them to dynamically query the knowledge base during inference without dataset-specific training.

1484 **Evaluation Metrics** We evaluate each KGQA system using 4 complementary Hit@1, Precision,
1485 Recall, and F1—under two answer-matching schemes: Exact Match (EM) and LASM. **EM** considers
1486 a prediction correct only if the model’s answer set exactly matches the ground-truth set after basic
1487 normalization, which includes lowercasing and alphabetically sorting answers. This is a strict
1488 string-level comparison that does not account for synonyms, paraphrases, or other forms of semantic
1489 equivalence. Hit@1 measures whether the model’s top-ranked answer appears in the ground-truth
1490 set. Precision, Recall, and F1 capture the degree of set overlap: Precision reflects the proportion of
1491 predicted answers that are correct, Recall captures the proportion of ground-truth answers that are
1492 retrieved, and F1 is their harmonic mean—together highlighting whether a model tends to over- or
1493 under-generate. **LASM** extends this evaluation by replacing literal comparison with a GPT-4o-mini
1494 verifier that determines whether the predicted and ground-truth answers are semantically equivalent.
1495 We then recompute all five metrics based on this semantic agreement. This two-tiered protocol offers
1496 a comprehensive view of model performance, balancing surface-level exactness with meaning-level
1497 correctness.

1498 **Cross-Dataset Evaluation.** To further contextualize model performance, we perform a cross-
1499 dataset comparison between our proposed KGQAGen-10k and two widely used KGQA benchmarks,
1500 WebQSP [73] and CWQ [58]. We evaluate four representative frameworks—RoG [33], GCR [34],
1501 ToG [57], and PoG [9]—under consistent inference settings. For WebQSP and CWQ, we report
1502 results from the respective original papers, while for KGQAGen-10k, we re-run each model using
1503 identical retrieval and reasoning configurations.

1504 As shown in Table 8, all models achieve high accuracy on traditional benchmarks (e.g., RoG reaches
1505 85.7 on WebQSP) but experience a substantial drop on KGQAGen-10k, where the same model
1506 attains only 21.3. This contrast highlights the increased reasoning complexity and stricter grounding
1507 requirements in our dataset. Moreover, the relative ranking of models remains consistent, but the
1508 gaps between them widen on KGQAGen-10k, indicating that it better differentiates model capabilities
1509 under more realistic and verifiable evaluation conditions.

Model	WebQSP (Hit@1)	CWQ (Hit@1)	KGQAGen-10k (Hit@1)
RoG	85.7	62.6	21.3
GCR	92.2	75.8	52.5
ToG	82.6	67.6	52.6
PoG	87.3	75.0	54.0

Table 8: Cross-dataset comparison of KGQA performance. Results on WebQSP and CWQ are taken from the respective original papers; KGQAGen-10k results are obtained under identical evaluation protocols.