# DIVKNOWQA: Assessing the Reasoning Ability of LLMs via Open-Domain Question Answering over Knowledge Base and Text

Anonymous ACL submission

#### Abstract

Large Language Models (LLMs) have exhibited impressive generation capabilities, but they suffer from hallucinations when solely 004 relying on their internal knowledge, especially when answering questions that require less commonly known information. Retrievalaugmented LLMs have emerged as a potential solution to ground LLMs in external knowledge. Nonetheless, recent approaches have primarily emphasized retrieval from unstructured text corpora, owing to its seamless integration into prompts. When using structured data such as knowledge graphs, most methods simplify it into natural text, neglecting the underlying structures. Moreover, a significant gap in the current landscape is the absence of a realistic benchmark for evaluating the effectiveness of grounding LLMs on heterogeneous knowledge sources (e.g., knowledge base and text). To fill this gap, we have curated a comprehensive dataset that poses two unique challenges: (1) **Two-hop multi-source questions** that require retrieving information from both opendomain structured and unstructured knowledge sources; retrieving information from structured knowledge sources is a critical component in correctly answering the questions. (2) Gener-027 ation of symbolic queries (e.g., SPARQL for Wikidata) is a key requirement, which adds another layer of challenge. Our dataset is created using a combination of automatic generation through predefined reasoning chains and human annotation. We also introduce a novel approach that leverages multiple retrieval tools, including text passage retrieval and symbolic language-assisted retrieval. Our model outperforms previous approaches by a significant margin, demonstrating its effectiveness in addressing the above-mentioned reasoning challenges.

# 1 Introduction

041

LLMs have shown exceptional performance in multi-hop question-answering (QA) tasks over text (TextQA) (Rajpurkar et al., 2018; Kwiatkowski



Figure 1: An example of the two-hop multi-source questions in DIVKNOWQA.

et al., 2019; Joshi et al., 2017; Trivedi et al., 2022a; Yang et al., 2018; Ho et al., 2020), tables (TableQA) (Yu et al., 2018; Zhong et al., 2017; Pasupat and Liang, 2015; Chen et al., 2019), and knowledge-bases (KBQA) (Gu et al., 2021; Yih et al., 2015; Talmor and Berant, 2018; Bao et al., 2016), where the supporting fact is contained in a single knowledge source – structured or unstructured. However, in many real-world scenarios, a QA system may need to retrieve information from both unstructured and structured knowledge sources; failing to do so results in insufficient information to address user queries.

While existing QA benchmarks provide diverse perspectives for evaluating models (Table 1), they are limited in assessing the performance of retrieval-augmented language models across heterogeneous knowledge sources in the following

Table 1: Comparing benchmarks for heterogeneous question-answering tasks. The column OpenR stands for open information retrieval, Human for human-written questions, EI for equitable importance of knowledge sources, and SGT for structured ground truth.

Dataset	KB	Text	Table	OpenR	Human	EI	SGT
HybridQA (Chen et al., 2021a)	X	~	~	X	1	~	X
OTT-QA (Chen et al., 2020)	X	1	1	1	~	~	×
NQ-Tables (Herzig et al., 2021)	×	1	1	X	×	×	×
TAT-QA (Zhu et al., 2021)	X	1	1	×	~	~	×
MultimodelQA (Talmor et al., 2021)	×	1	1	X	1	×	×
Manymodelga (Hannan et al., 2020)	X	1	1	×	~	×	×
FinQA (Chen et al., 2021b)	×	1	1	X	1	1	×
HetpQA (Shen et al., 2022)	×	1	1	X	1	×	×
CompMix (Christmann et al., 2023)	X	1	1	1	~	×	×
WikiMovies-10K (Miller et al., 2016)	1	1	×	1	1	×	×
MetaQA (Zhang et al., 2018)	1	1	×	1	1	×	×
DIVKNOWQA (Ours)	1	1	×	1	1	1	1

aspects: (1) Closed-book QA: Closed-book questions do not accurately reflect the real-world setting where individuals generally have access to diverse knowledge sources on the Internet; (2) Automatically generated data: The lack of human verification results in erroneous data; (3) Imbalanced emphasis across different knowledge sources: Current benchmarks feature knowledge sources with varying levels of importance. Answers may be found in multiple sources, leading models to prioritize textual sources while underutilizing structured knowledge sources; (4) Suboptimal use of structured knowledge: Structured knowledge sources are typically treated as textual sources by linearizing triplets from the knowledge base or rows/columns from tables, missing the opportunity to fully realize the benefit of highly-precise structured knowledge by probing them via symbolic queries.

063

071

081

089

094

100

101

Despite the inherent challenges, being able to generate structured queries effectively can offer a number of benefits. First, unlike a query to retrieve text passages, the structured query itself can share the responsibility of reasoning (Liu et al., 2022). For example, for the question "How many awards has Neil Armstrong received?", to get an answer from a knowledge base such as Wikidata (Vrandečić and Krötzsch, 2014), a SPARQL query (Pérez et al., 2006) can use an aggregation function to return the numerical number as the final result as shown in Figure 1. In contrast, a text retriever needs to locate all the relevant passages and rely on a reader module to get the final result. The commonly used readers often come with an input length constraint. The number of returned passages could be too many to fit into the reader's context, causing a wrong answer. Even when the context length is not an issue, even the best LLMs have difficulties in locating the answer (Liu et al., 2023a). Besides, there is less room for ambiguity in structured queries. For example, a dense retriever cannot easily distinguish between similar song titles such as "I'll be good to you" and "I have been good to you" by different singers. On the other hand, given the right identifier of the entity, the structured knowledge search engine can return the relevant information for the exact entity. 102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

In this work, we propose DIVKNOWQA, a novel fact-centric multi-hop QA benchmark that requires models to utilize heterogeneous knowledge sources equitably in order to answer a question. We perform the first study to assess the reasoning ability of LLMs, via jointly exploiting open-domain QA over heterogeneous knowledge sources. In particular, we have chosen Knowledge Base (KB) as our primary case study for the structured source, and we have created a dataset comprising 940 human-annotated examples. Additionally, each entry in our dataset includes a corresponding symbolic SPARQL query to facilitate the retrieval of information from the KB. To generate the questions, we construct a question collection pipeline comprising three key steps: text-based QA sampling, KB question generation, and question composition, all while minimizing the need for human annotation efforts.

To set up a baseline, in addition to benchmarking on standard and tool-augmented LLMs, we propose a Diverse rEtrieval Tool Augmented LLM (DETLLM) to address the challenges posed by DIVKNOWQA. DETLLM decomposes a multihop question into multiple single-hop questions, and adopts two novel strategies: (1) *symbolic query generation* to retrieve supportive text from a KB by transforming a single-hop natural question into a SPARQL query, and (2) *retrieval tool design*, which includes a textual retriever and a symbolic query generation tool to recall relevant evidence from heterogeneous knowledge sources. Our method shows improvements of up to 4.2% when compared to existing methods.

# 2 The DIVKNOWQA Dataset

#### 2.1 Dataset Collection

Our goal is to create a method for generating complex questions from diverse knowledge sources,145plex questions from diverse knowledge sources,146making each source indispensable; and we aim147to do so with a minimal human annotation effort.148Additionally, we wish to provide Wikidata entity149and relation IDs to support structured query-based150knowledge retrieval. Due to the page limit, Figure151

3 from Appendix A.5 depicts our proposed method. 152 We first sample a single-hop text question from 153 the Natural Question dataset (Kwiatkowski et al., 154 2019) as an anchor, to which we link to a relevant 155 Wikidata triplet. Then single-hop KB questions are 156 generated based on the sampled triplets thereby 157 using the anchor question to automatically com-158 pose a heterogeneous multi-hop question. Human 159 annotators finally verify the quality of the machine-160 generated question and rewrite the question that 161 needs revision. In the following, we elaborate on 163 the steps.

Natural Questions as Anchors The Natural Question (NQ) dataset is a question-answering 165 dataset containing tuples of (question, 166 answer, title, passage), where title 167 and passage are respectively the title of the Wikipedia page and the passage containing the 169 answer. The questions in NQ were collected from 170 real-world user queries issued to the Google search engine, and it contains 307K training examples. 172 We concentrate on constructing a multi-hop dataset 173 linked through the initial step's single-hop answer. 174 To achieve this, we extract question-answer pairs 175 where the question contains a succinct answer of 176 177 up to 5 words to ensure the quality of the resulting composed question. 178

Linking Natural Questions to Wikidata We 179 adopt the notion of bridge entity from Yang et al. (2018) to describe the single-hop answer in the 181 initial step when breaking down a multi-hop 182 question. We explore two linking options, each involving a unique choice of bridging an entity to connect the natural question to Wikidata. We explain the options using the example question "Who plays Mary Poppins in Mary Poppins Returns?" with the answer "Emily Blunt". (a) Text  $\rightarrow$  KB Approach: We 189 treat the answer "Emily Blunt" as the bridge 190 entity, and search for a Wikidata triplet where 191 "Emily Blunt" is the subject, for example, (Emily Blunt, sibling, Felicity 193 Blunt). (b) KB  $\rightarrow$  Text Approach: In this 194 alternative method, we recognize the question 195 entity that exists in Wikidata, in this case, "Mary Poppins Returns", as the bridge entity. For simplicity, we only consider the entity mentioned 198 in the Wikipedia title. We then link to the 199 Wikidata triplet using it as the *object*, leading to triplets such as "(William Weatherall

Wilkins, present in work, Mary Poppins Returns)". 202

203

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

Selection of KB Triplets To maintain an equal 204 emphasis on both structured and unstructured 205 knowledge sources, we implement a meticulous 206 selection process for KB triplets to ensure that the 207 associated knowledge cannot be easily obtained 208 by merely retrieving information from the textual 209 source (Wikipedia passages). We retain triplets 210 (sub, relation, obj), where either the subject sub 211 is not linked to a Wikipedia page or the object *obj* 212 does not exist within the Wikipedia page associated 213 with the subject. This ensures that simply retriev-214 ing the Wikipedia passage for the *sub* is unlikely 215 to yield an answer to a question involving sub and 216 *obj*, thereby requiring the model to utilize the KB. 217 Furthermore, when generating questions in the KB 218  $\rightarrow$  Text linking option, we selectively retain triplets 219 where only one object is associated with the given relation and subject This approach ensures the com-221 pleteness and uniqueness of the reasoning chain. 222 For example, given a composed question, "Who 223 plays Mary Poppins in Lin-Manuel Miranda's no-224 table work?", "Mary Poppins Returns" is one of 225 the notable works from "Lin-Manuel Miranda". By 226 querying KB given the subject "Mary Poppins Re-227 turns" and the relation "notable work", we will lo-228 cate multiple answers rather than the single bridge 229 entity "Mary Poppins Returns", posing a challenge 230 to infer the second sing-hop question "Who plays 231 Mary Poppins in Mary Poppins Returns?".

Generating Single-Hop KB Questions We then create single-hop questions from the selected (*sub*, *relation*, *obj*) triplets. These questions are designed to emphasize the relationship between *sub* and *obj*, with *obj* being the expected answer. For instance, for the KB triplet "(Emily Blunt, sibling, Felicity Blunt)", we expect to generate a question like "Who is the sibling of Emily Blunt?". For this, we employ the gpt-turbo-3.5 LLM from OpenAI; the prompt can be found in Appendix A.1.

Generating Heterogeneous Multi-Hop Questions In this stage, we wish to create a multi-hop question by composing a textual question and a KB question. We generate such *heterogeneous* questions by carefully chaining two single-hop questions together. DIVKNOWQA supports three question types: short entity, yes/no, and aggregate questions, and two question composition orders: Text

 $\rightarrow$  KB and KB  $\rightarrow$  Text. This combination results in a total of five question types, as we construct aggre-253 gate questions following only the Text  $\rightarrow$  KB order. 254 We employ gpt-turbo-3.5 as a question composer to connect two single-hop questions. This is achieved by substituting the entity mentioned in the 257 outer question with a rephrased version of the first question. The prompt for generating the multi-hop questions is given in Appendix A.2. Our generation method for different question types is discussed as 261 follows. 262

**Short Entity Question** We use a factoid entity as the final answer. The final answer can be the object from Wikidata or the factoid answer from NQ.

263

264

Yes/No Question In contrast to Short Entity questions, Yes/No questions involve an additional step. Initially, the original question is reformulated into a verification-style question typically starting with 269 phrases like "Is/Was/Were/Does/Do/Did". This new question includes a candidate answer for 271 verification purposes. For instance, let's consider 272 the original question "What grade were they in High School Musical 1?" with a known answer of "juniors". To create a verification question, we might rephrase it as 276 "Were they seniors in High School 277 Musical 1?" and include the verifying answer 278 "seniors" within the question. Generating the candidate answer for verification can be a complex task as it requires choosing a verifying 281 answer that aligns well with the context of the question. Sampling incorrect distractors as verifying answers is also a part of the process. These distractors should be incorrect but closely related to the answer, and they are generated by prompting qpt-turbo-3.5. This approach ensures that the verification process is robust and accurate, preventing situations where the verifying answer deviates from the question's context and 290 potentially leads to a simplistic answer "no" during evaluation.

293Aggregate QuestionWe formulate aggregating294questions in the "Text  $\rightarrow$  KB" composition295order, where the outermost question pertains to296counting the number of associated triplets based297on the given subject and relation. For instance,298the outermost question "How many awards299does Milton Friedman receive?"300arises from the KB triplets of the form (Milton301Friedman, award received, award

name)" with 10 such award name objects. In such cases, we leverage the aggregate feature offered by the structural query (i.e., SPARQL).

302

303

304

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

331

332

333

334

335

336

337

338

339

341

342

343

344

345

346

347

# 2.2 Human Annotation

We recruited five individuals, three undergraduate and two graduate students with experience in the field of NLP for data verification and annotation. Each question underwent a verification and rewriting process involving two annotators. To mitigate any potential annotation bias, we presented each question to both annotators, with the order of examples shown to annotators randomized. Annotators were tasked with assessing the quality of KB question generation, and they had three options to choose from: "Accept", "Revise", or "Reject". When a question required revision, annotators were instructed to make modifications while preserving the focus on the subject and relation and keeping the answer unchanged. Additionally, they were responsible for evaluating the quality of complex questions and providing necessary revisions. The instruction provided to human annotators is shown in Appendix A.4. Annotators were duly compensated for their valuable contributions to our study. Out of 1,000 examples that were annotated, 757 examples received unanimous approval, 183 underwent revisions, and 60 were rejected. Both unanimously accepted and revised examples were included in the dataset.

### 2.3 Dataset Statistics and Analysis

In this section, we analyze the question types and KB single-hop relation types in DIVKNOWQA.

**Question Central Word** Taking inspiration from (Yang et al., 2018), we designate the first three words of a question as the Center Question Words (CQW). We adopt this approach because our questions typically do not contain comparison queries, and a majority of question words are found at the beginning of the question. Due to the page limit, Figure 4(a) in Appendix A.6 provides a visual representation of CQW in DIVKNOWQA.

**KB Relation Types** We also analyze the distribution of relations by counting the frequency of different relations that appear in the KB triplets used to construct the single-hop KB questions. Due to the page limit, Figure 4(b) in Appendix A.6 features the distribution of diverse relations.

A How many awards has the first person to walk on the moon received?	Answer: 26
Ú LUI	M
Rationale: Decompose the question to answer the following single-hop questions. 1. who was the first person to walk on the moon? 2. how many awards has this person received?	Rationale: The first person to walk on the moon is Neil Armstrong. The second step is to answer how many awards he has received.
Query: Who was the first person to walk on the moon?	Query: How many awards has Neil Armstrong received?
Paragraphs [1] This was accomplished with two US pilot- astronauts flying a Lunar Module on each of six NASA missions across a 41-month period starting on 20 July 1969 UTC, with <b>Neil</b> Armstrong and [2] <b>Neil Armstrong</b> (the first person to walk on the moon) Linearized Triplets [3] <b>Neil Armstrong</b> ; Description; Armstrong was an American astronaut and the first person to walk on the Moon [4] John Young; Spaceflight Astronaut Missions; Apollo 16 SPARQL None	Paragraph         [1] She was the first female newscaster on television in Los Angeles and the West Coast She has received many awards and honors         [2] Apollo 11 was the spaceflight that landed the first two people on the Moon. Commander Neil Armstrong and Lunar Module         Linearized Triplets         [3] Neil Armstrong; Award Received Silver Buffalo Award. Neil Armstrong; Award Received; Livingstone Medal         [4] Neil Armstrong; Award Received; Air Medal . Neil Armstrong; Award Received; Collier Trophy         SPARQL         SELECT (COUNT(?award) as ?count)         WHERE { wd:Q1615 wdt:P166 ?award. }
Rationale: Neil Armstrong was the first person to walk on the moon and he has received 26 awards.	Answer: 26

Figure 2: The illustration of DETLLM to instruct LLMs for addressing multi-source multi-hop questions.

Anecdotal Examples for Representative Types 349 Due to the length limit, in Appendix A.7, Table 5 350 presents illustrative examples drawn from the DI-351 VKNOWQA benchmark for each of our five question composition types. These examples serve to showcase how our dataset necessitates information 354 retrieval from diverse sources in varying orders. Additionally, they highlight that the answer types 357 require models to perform tasks such as answer span extraction, candidate answer verification, and information aggregation based on relevance.

# 3 DETLLM: Diverse Retrieval Tool Augmented LLM

361

362

363

364

367

369

We now introduce our diverse retrieval tool augmented LLM (DETLLM) and show its promising capability on the proposed DIVKNOWQA benchmark by unifying the retrieval ability from the structured and unstructured knowledge sources.

To tackle a complex question, we follow the chain-of-thought (CoT) framework (Wei et al., 2022) to decompose a complex question into singlehop questions where each single-hop question is knowledge-intensive, requiring supportive fact retrieval from a knowledge source.

We design a retrieval tool capable of retrieving from heterogeneous knowledge sources. For unstructured text knowledge, a dense passage retriever (Izacard et al., 2022) is employed to retrieve relevant passages. For structured knowledge, we consider two modalities of structured knowledge to maximize the relevant information coverage. First, we transform the structured data into text passages by linearizing the relation triplets into passages in which case a sparse text retriever can be used to detect similar sources. Second, we propose a symbolic query generation module to map a natural language query to a structured query (e.g., SPARQL) to directly query against the KB (e.g., Wikidata). The benefits are twofold: (1) pinpointing precise knowledge, and (2) leveraging the compositionality of the query language and reducing the mere reliance on the language model's reasoning responsibility. Figure 2 shows the DETLLM flow for querying an LLM.

373

374

375

376

377

378

380

381

382

383

384

385

386

387

388

389

390

391

### 3.1 Question Decomposition and Planning

394

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425 426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

Our approach to answering a complex multi-hop question is inspired by the conceptual framework of DSP (Khattab et al., 2022). When dealing with a question that involves n hops, we query the LLM *n* times to generate retrieval queries and retrieve information from a knowledge source. The final query is used to ultimately arrive at the final answer, utilizing the retrieved passage to answer the last single-hop question. This process results in a total of n + 1 interactions with the LLM. At the *j*-th LLM prompting, the LLM's task is to utilize the previously retrieved information to answer the j-1 single-hop question. It then dissects the original question Q into the *j*-th subsequent single-hop questions  $q_i$ , which serve as a retriever query to gather information from a knowledge source.

#### 3.2 Multi-Source Knowledge Retrieval

In addressing the single-hop questions, our approach entails searching across diverse knowledge sources to gather supporting facts. To answer the subsequent single-hop question  $q_j$ , we begin by having the LLM generate semantically diverse queries, denoted as  $Query_j = \{query_1^j, \ldots, query_t^j\}$ . We set the LLM decoding temperature to 0.7 to sample diverse queries.

In our approach, we treat unstructured and structured knowledge separately and retrieve relevant information from both knowledge sources. As mentioned, for unstructured knowledge, we use a dense retriever Contriever (Izacard et al., 2022) to retrieve relevant passages, while for structured knowledge, we retrieve relevant information from both textual and structured formats. The preparation of the textual knowledge base involves linearizing KB triplets (*sub*, *relation*, *obj*) into a string format "sub relation obj" after which we create a retrieval index for efficient passage retrieval using a sparse text retriever BM25 (Robertson et al., 2009). The Contriever, trained on natural language corpus, is adaptable to unstructured knowledge but struggles when faced with linearized structured knowledge because it lacks natural language formatting. In contrast, the sparse retriever BM25 performs better with structured knowledge by using a keywordbased search methodology. We show the ablation study in Section 4.3.

In addition, we generate SPARQL queries to execute against the Wikidata engine to retrieve further relevant information. Our retrieval tool thus com-

Table 2: Answer and Sub-Step Retrieval Accuracy onDIVKNOWQA.

	EM	F1	Recall	H1-R	H2-R
Vanilla Prompt	26.0	28.3	26.8	42.2	-
ReAct	16.1	18.4	19.0	-	-
DSP	27.9	31.0	31.2	57.6	41.2
DETLLM (our)	32.1	35.7	35.6	70.1	47.1

prises three components: a sparse retriever, a dense retriever, and a symbolic query language generation module. These elements collectively enable the comprehensive retrieval of information from heterogeneous knowledge sources. 443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

### 3.3 Multi-Source Knowledge Ranking

To consolidate the retrieved information obtained from the tool, we perform a ranking of information from various knowledge sources. The goal is to select the top-k most relevant pieces of information. This selection is necessary because of the inherent length constraint of the language model, which prevents us from incorporating all the retrieved information into the prompt. To achieve this ranking, our approach leverages the off-of-shelf crossencoder model (Reimers and Gurevych, 2019) to assess the relevance of each piece of retrieved information in the context of a single-hop question. We use the sentence-transformers package implementation with the model checkpoint cross-encoder/ms-marco-MiniLM-L-6  $-v2.^{1}$ 

# 4 Benchmarking

## 4.1 Experimental Setup

**Baselines** (1)ChatGPT (OpenAI, 2023): We employ OpenAI's ChatGPT model (qpt-3.5-turbo) by single-step query inputting the question and retrieved-context and obtaining its response as the final answer. (2) DSP (Khattab et al., 2022): We apply the demonstrate-search-predict framework to iteratively address complex QA tasks with the assistance of retrieved context. (3) ReAct (Yao et al., 2023): It leverages a synergistic approach, combining reasoning with tool usage. It involves verbally generating a reasoning trace and issuing the necessary commands to invoke a tool, which then takes action accordingly. We use gpt-3.5-turbo as the backbone model.

**Evaluation Metrics** To assess the accuracy and 482 relevance of various models for factoid questions, 483 we rely on established metrics. We report the exact 484 match and F1 score for final answer quality, follow-485 ing the methodology of (Yang et al., 2018). Besides, 486 we report the Recall score indicating whether the 487 ground-truth answer is a substring in prediction 488 since LLM may generate extra information. In ad-489 dition, we report the retrieval accuracy for each 490 decomposed single-hop question denoted as H1-R 491 and H1-2 for the two-hop question. 492

> **Implementation Details** To ensure a comprehensive and equitable comparison, we offer baseline model access to both structured knowledge and unstructured knowledge as retrieval sources. In the case of the baseline model, the KB is converted into linearized passages, which are then combined with the unstructured knowledge, creating a unified source for retrieval. We use BM25 (Robertson et al., 2009) and Contriever (Izacard et al., 2022) as sparse and dense retrieval tools respectively. Unless specified otherwise, we experiment with a few-shot prompt that includes three humanannotated demonstrations along with task instructions to guide the model generation process.

## 4.2 Main Results

493

494

495

496

497

498

499

504

506

508

510

511

512

514

515

517

518

519

522

524

525

526

527

Comparing with State-of-the-Art LLMs Table 2 presents the model performance results on the DIvKNOWQA. ReAct exhibits lower performance compared to the Vanilla prompt. The retrieval tool created for ReAct is specialized for querying unstructured knowledge. As the presence of irrelevant passages distracts the LLM (Chen et al., 2023; Mallen et al., 2023), the iterative reasoning accumulates errors, leading to less accurate answers. Conversely, DSP outperforms both Vanilla Prompt and ReAct, thanks to its robust search module designed to engage with frozen retrievers. DSP enhances a single retrieval query into multiple queries, employing a fusion function to rank candidate passages and identify the most relevant one. However, the search module cannot effectively retrieve structured knowledge. Our model stands out as the topperforming model, demonstrating its capability to generate symbolic language for retrieval from structured knowledge.

528Retrieval PerformanceTable 2 also presents the529single-step retrieval accuracy. Among the base-530line methods, comparing single-step generation

Table 3: Ablation study on the retrieval strategy.

	EM	F1	Recall	H1-R	H2-R
w/o SPARQL					
Text-KB(Sparse)	27.9	31.0	31.2	57.6	41.2
Text-KB(Dense)	22.7	26.1	26.9	54.9	32.0
Text(Sparse)-KB(Sparse)	26.4	29.8	30.2	60.0	41.2
Text(Dense)-KB(Sparse)	30.7	35.0	35.5	68.9	46.8
w/ SPARQL					
Text-KB(Sparse)	28.8	31.9	32.9	58.0	42.9
Text-KB(Dense)	31.2	34.7	35.9	64.3	42.6
Text(Sparse)-KB(Sparse)	28.5	31.8	32.0	61.5	42.1
DETLLM (our)	32.1	35.7	35.6	70.1	47.1

e.g. Vanilla Prompt with the multi-step generation e.g. ReAct and DSP, the retrieval accuracy increases due to the decomposed query from the multi-step generation process. On the other hand, the DETLLM shows stronger retrieval performance compared to DSP due to the careful retrieval tool design, the unstructured and structured knowledge is treated separately. This finding underscores the importance of having a robust retrieval strategy to provide reliable and focused information, grounding the LLM on relevant supportive facts.

### 4.3 Discussion

Ablation Study Table 3 presents the results of an ablation study involving three key factors: a) the integration of heterogeneous knowledge sources, b) the choice between dense and sparse retrievers, and c) the incorporation of SPARQL. Our findings indicate that optimal performance is achieved when handling heterogeneous knowledge sources separately, combined with careful retriever tool selection. The unsupervised dense retriever (i.e., Contriever), trained on natural language corpus, demonstrates adaptability to unstructured knowledge but loses its advantage when dealing with linearized structured knowledge due to the absence of natural language formatting. Conversely, the sparse retriever BM25 performs better on structured knowledge, relying on keyword-based search methodologies. Furthermore, the SPARQL tool consistently outperforms its counterparts in all settings, showcasing improvements regardless of the integration of knowledge sources and the choice of retriever.

Table 4: Comparison between the closed book settingand open domain retrieval.

	EM	F1	Recall	H1-R	H2-R
Closed Book	30.2	33.8	31.2	-	-
DetLLM	32.1	35.7	35.6	70.1	47.1

561

562

**Comparing with Closed-book LLM** Table 4 563 presents a comparison between DETLLM and 564 LLM performance in the closed-book setting, 565 where no external knowledge is accessible. We demonstrate that DETLLM exhibits improvements in scenarios distinct from the closed-book setting. 568 We observe that only 50.8% of examples answered 569 correctly by our DETLLM are also present in the closed-book setting, highlighting the orthogonal 571 performance of DETLLM compared to the closed-572 book setting. The combination of correctly an-573 swered examples accounts for 45.4% of the entire 574 dataset. One plausible hypothesis is that the closed 575 book setting enables the LLM to access knowledge 576 stored in its memory, reducing the impact of re-577 triever errors. We also suggest a potential research direction, which involves designing a strategy to switch between the closed book setting and open domain retrieval to achieve optimal performance. 581

> Additional Discussion Due to the page length limit, we put additional discussion in the Appendix A.8. We present the SPARQL generation analysis in Table 6 and oracle retrieval performance in Table 7.

# 5 Related Work

584

585

590

591

593

596

599

603

604

610

611

# 5.1 Assessing the Reasoning Ability of LLMs

LLMs (Brown et al., 2020; Touvron et al., 2023; Nijkamp et al., 2023) have exhibited notable advancements in their capabilities, particularly in the domain of reasoning skills. These skills encompass various categories, including inductive reasoning (Wang et al., 2023; Yang et al., 2022), deductive reasoning (Creswell et al., 2023; Han et al., 2022), and abductive reasoning (Wiegreffe et al., 2022; Lampinen et al., 2022), depending on the type of reasoning involved. Current research efforts have predominantly focused on evaluating LLMs in the context of open-ended multi-hop deductive reasoning. These scenarios involve complex question-answering tasks (Yang et al., 2018; Gu et al., 2021; Trivedi et al., 2022b; Liu et al., 2023b) and fact-checking (Jiang et al., 2020). Notably, our work contributes to this landscape by introducing an additional layer of complexity: the integration of multi-hop and multi-source reasoning. In our approach, we retrieve supporting facts from heterogeneous knowledge sources, further enhancing the challenges posed to LLMs in this deductive reasoning context.

### 5.2 Retrieval-Augmented LLMs

Retrieval-Augmented Large Language Models (RALLMs) are semi-parametric models that integrate both model parameters and a non-parametric datastore to make predictions. RALLMs enhance LLMs by updating their knowledge (Izacard et al., 2023; Khandelwal et al., 2020; Yavuz et al., 2022; Mallen et al., 2023), providing citations to support trustworthy conclusions (Menick et al., 2022; Gao et al., 2023). RALLMs can retrieve information in an end-to-end fashion within a latent space (Khandelwal et al., 2020, 2021; Min et al., 2023), or they can follow the retrieve-then-read paradigm, leveraging an external retriever to extract information from textual sources (Ram et al., 2023; Khattab et al., 2022). Our approach adheres to the retrievethen-read paradigm, with a specific emphasis on multi-source retrieval, advocating for structured knowledge retrieval through symbolic generation.

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

## 6 Conclusion

We introduce the DIVKNOWQA, designed to evaluate the proficiency of question-answering systems, especially those enhanced by retrieval tools, in addressing knowledge-intensive questions with a strong emphasis on multi-hop multi-source retrieval. This dataset is constructed through automated data generation and subsequent human verification, minimizing manual effort. Our evaluation encompasses both standard LLMs and LLMs augmented with retrieval tools. Notably, we identify that this task presents a new challenge for state-of-the-art models due to the demand for structured knowledge retrieval and the inherent lack of prior knowledge in this context. To tackle this challenge, we propose the DETLLM, which incorporates diverse retrieval tools including innovative symbolic query generation for retrieving information from the structured knowledge source. In the future, we are keen on enhancing LLMs' capabilities in understanding and generating symbolic language, as well as exploring methods to improve performance on knowledge-intensive and complex question-answering tasks.

## Limitations

One limitation of our proposed DETLLM is that the retrieval tool is used in each decomposed single-hop question-answering step. A further step involves investigating when the large language model truly requires retrieval knowledge, rather

768

than invoking the tool at every step. Recent research (Mallen et al., 2023) has indicated that LLMs derive substantial benefits from general domain knowledge but may encounter challenges when dealing with long-tail knowledge because LLMs' memorization is often limited to popular knowledge. Future work can address the issue of uncertainty in LLMs' reliance on retrieval tools, aiming to optimize tool usage efficiently and establish trustworthiness in the process.

Another limitation is the need to explore the impact of extended prompts on retrieval-augmented language models. Recent research has revealed that LLMs can be susceptible to recency bias (Liu et al., 2023a). Furthermore, a study (Peysakhovich and Lerer, 2023) indicates that documents containing the ground truth answer tend to receive higher attention, suggesting that reordering documents by placing the highest-attention document at the forefront can enhance performance. Thus, an avenue for further investigation of whether document reordering strategies, based on attention mechanisms, can be employed to improve retrieval-augmented LM performance on multi-source multi-hop QA task.

#### References

661

667

670

672

674

675

679

684

702

703

704

706

707

712

- Junwei Bao, Nan Duan, Zhao Yan, Ming Zhou, and Tiejun Zhao. 2016. Constraint-based question answering with knowledge graph. In *Proceedings of COLING 2016, the 26th international conference on computational linguistics: technical papers*, pages 2503–2514.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Wenhu Chen, Ming-Wei Chang, Eva Schlinger, William Wang, and William W Cohen. 2021a. Open question answering over tables and text. *Proceedings of ICLR 2021*.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. 2019. Tabfact: A large-scale dataset for table-based fact verification. *arXiv preprint arXiv:1909.02164*.
- Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang. 2020. HybridQA: A dataset of multi-hop question answering over tabular and textual data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*,

pages 1026–1036, Online. Association for Computational Linguistics.

- Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Chatcot: Tool-augmented chain-of-thought reasoning on chatbased large language models.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021b. FinQA: A dataset of numerical reasoning over financial data. In *Proceedings* of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 3697–3711, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Philipp Christmann, Rishiraj Saha Roy, and Gerhard Weikum. 2023. Compmix: A benchmark for heterogeneous question answering. *arXiv preprint arXiv:2306.12235*.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 2023. Selection-inference: Exploiting large language models for interpretable logical reasoning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5,* 2023. OpenReview.net.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling large language models to generate text with citations. *arXiv preprint arXiv:2305.14627*.
- Yu Gu, Sue Kase, Michelle Vanni, Brian Sadler, Percy Liang, Xifeng Yan, and Yu Su. 2021. Beyond iid: three levels of generalization for question answering on knowledge bases. In *Proceedings of the Web Conference 2021*, pages 3477–3488.
- Simon Jerome Han, Keith Ransom, Andrew Perfors, and Charles Kemp. 2022. Human-like property induction is a challenge for large language models. In *Proceedings of the Annual Meeting of the Cognitive Science Society*.
- Darryl Hannan, Akshay Jain, and Mohit Bansal. 2020. Manymodalqa: Modality disambiguation and qa over diverse inputs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7879–7886.
- Jonathan Herzig, Thomas Müller, Syrine Krichene, and Julian Eisenschlos. 2021. Open domain question answering over tables via dense retrieval. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 512–519, Online. Association for Computational Linguistics.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multihop QA dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*,

- 769 770 771 774 776 777 778 779 780 781 783 787 788 789 790 791 793 794 800 810

824

811

812

813 814

815

817

816

ton Lee, et al. 2019. Natural questions: a benchmark

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-

predict: Composing retrieval and language models for knowledge-intensive nlp. arXiv preprint arXiv:2212.14024.

466.

on Learning Representations (ICLR). Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through Memorization: Nearest Neighbor Language Models. In International Conference on Learning Representations (ICLR).

bor machine translation. In International Conference

Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts,

and Matei Zaharia. 2022. Demonstrate-search-

for question answering research. Transactions of the

Association for Computational Linguistics, 7:453-

Andrew Lampinen, Ishita Dasgupta, Stephanie Chan,

Kory Mathewson, Mh Tessler, Antonia Creswell,

James McClelland, Jane Wang, and Felix Hill. 2022.

Can language models learn from explanations in con-

text? In Findings of the Association for Computa-

tional Linguistics: EMNLP 2022, pages 537-563,

Abu Dhabi, United Arab Emirates. Association for

Computational Linguistics.

Canada. Association for Computational Linguistics. Zettlemoyer, and Mike Lewis. 2021. Nearest neigh-

ume 1: Long Papers), pages 1601–1611, Vancouver, Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke

pages 6609-6625, Barcelona, Spain (Online). Inter-

national Committee on Computational Linguistics.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebas-

tian Riedel, Piotr Bojanowski, Armand Joulin, and

Edouard Grave. 2022. Unsupervised dense informa-

tion retrieval with contrastive learning. Transactions

augmented language models. Journal of Machine

Yichen Jiang, Shikha Bordia, Zheng Zhong, Charles

Dognin, Maneesh Singh, and Mohit Bansal. 2020.

HoVer: A dataset for many-hop fact extraction and

on Machine Learning Research.

Learning Research, 24(251):1-43.

Linguistics.

the Association for Computational Linguistics (Vol-

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In Proceedings of the 55th Annual Meeting of

claim verification. In Findings of the Association for Computational Linguistics: EMNLP 2020, pages 3441-3460, Online. Association for Computational

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu. Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval

- Ye Liu, Semih Yavuz, Rui Meng, Dragomir Radev, Caiming Xiong, and Yingbo Zhou. 2023b. Answering complex questions over text by hybrid question parsing and execution. arXiv preprint arXiv:2305.07789.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 9802-9822, Toronto, Canada. Association for Computational Linguistics.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, et al. 2022. Teaching language models to support answers with verified quotes. arXiv preprint arXiv:2203.11147.

Alexander Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. 2016. Key-value memory networks for directly reading documents. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1400–1409, Austin, Texas. Association for Computational Linguistics.

Sewon Min, Weijia Shi, Mike Lewis, Xilun Chen, Wentau Yih, Hannaneh Hajishirzi, and Luke Zettlemover. 2023. Nonparametric masked language modeling. In Findings of the Association for Computational Linguistics: ACL 2023, pages 2097–2118, Toronto, Canada. Association for Computational Linguistics.

Erik Nijkamp, Tian Xie, Hiroaki Hayashi, Bo Pang, Congying Xia, Chen Xing, Jesse Vig, Semih Yavuz, Philippe Laban, Ben Krause, et al. 2023. Xgen-7b technical report. arXiv preprint arXiv:2309.03450.

OpenAI. 2023. Introducing chatgpt.

Panupong Pasupat and Percy Liang. 2015. Compositional semantic parsing on semi-structured tables. arXiv preprint arXiv:1508.00305.

Jorge Pérez, Marcelo Arenas, and Claudio Gutierrez. 2006. Semantics and complexity of sparql. In International semantic web conference, pages 30-43. Springer.

Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023a. Lost in the middle: How language models use long contexts. arXiv preprint arXiv:2307.03172.

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

Ye Liu, Semih Yavuz, Rui Meng, Dragomir Radev, Caiming Xiong, and Yingbo Zhou. 2022. Uni-parser: Unified semantic parser for question answering on knowledge base and database. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 8858–8869.

991

935

- Alexander Peysakhovich and Adam Lerer. 2023. Attention sorting combats recency bias in long context language models. *arXiv preprint arXiv:2310.01427*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*.

891

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

924

925

926

927

930

931

- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389.
- Xiaoyu Shen, Gianni Barlacchi, Marco Del Tredici, Weiwei Cheng, Bill Byrne, and Adrià Gispert.
  2022. Product answer generation from heterogeneous sources: A new benchmark and best practices. In *Proceedings of the Fifth Workshop on e-Commerce* and NLP (ECNLP 5), pages 99–110, Dublin, Ireland. Association for Computational Linguistics.
- Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. *arXiv preprint arXiv:1803.06643*.
- Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. 2021. Multimodal{qa}: complex question answering over text, tables and images. In *International Conference on Learning Representations*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022a. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.

- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022b. MuSiQue: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.
- Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85.
- Ruocheng Wang, Eric Zelikman, Gabriel Poesia, Yewen Pu, Nick Haber, and Noah D Goodman. 2023. Hypothesis search: Inductive reasoning with language models. *arXiv preprint arXiv:2309.05660*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark Riedl, and Yejin Choi. 2022. Reframing human-AI collaboration for generating free-text explanations. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 632–658, Seattle, United States. Association for Computational Linguistics.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2022. Language models as inductive reasoners. *arXiv preprint arXiv:2212.10923*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net.
- Semih Yavuz, Kazuma Hashimoto, Yingbo Zhou, Nitish Shirish Keskar, and Caiming Xiong. 2022. Modeling multi-hop question answering as single sequence prediction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics.
- Scott Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. 2015. Semantic parsing via staged query graph generation: Question answering with knowledge base. In *Proceedings of the Joint Conference of the 53rd Annual Meeting of the ACL and the 7th International Joint Conference on Natural Language Processing of the AFNLP*.

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. *arXiv preprint arXiv:1809.08887*.

992

993

994

995

998

999

1000

1001

1002

1003

1005

1006

1007 1008

1009

1010

1011

1012

1013

1014

1015

1016

- Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander Smola, and Le Song. 2018. Variational reasoning for question answering with knowledge graph. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv preprint arXiv:1709.00103*.
- Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3277–3287, Online. Association for Computational Linguistics.

# A Appendix

# A.1 Single-Hop Knowledge Base Question 1018 Generation Prompt 1019

1017

1047

1048

Prompt 1: Single-Hop Knowledge Base Question Generation

Instruction: Question generation given	1020
the following information:	1021
1) Answer	1022
2) Short relation between the question	1023
entity and the answer	1024
3) Question entity.	1025
	1026
IMPORTANT: The answer must be avoided	1027
in the question.	1028
	1029
Answer: Jacques Boigelot;	1030
Relation: director;	1031
Question Entity: Peace in the Fields;	1032
Question: Who directs Peace in the	1033
Fields?	1034
	1035
Answer: Academy Award for Best Sound	1036
Mixing;	1037
Relation: award received;	1038
Question Entity: Douglas Shearer;	1039
Question: Which award does Douglas	1040
Shearer receive?	1041
	1042
Answer: Rio de Janeiro;	1043
Relation: place of birth;	1044
Question Entity: David Resnick;	1045
Question: Where was David Resnick born?	1046

# A.2 Multi-Hop Complex Question Generation Prompt

Prompt 2: Multi-Hop Complex Question Generation

Instruction: Compose 2 single-hop	1049
questions into a 2-hop question	1050
given:	1051
1) Hop1 question	1052
2) Hop1 answer	1053
3) Hop2 question	1054
-)F- 1	1055
Hopl question: Who said a rose by any	1056
other name would smell just as	1057
sweet?	1058
Hopl answer: Juliet	1059
Hop2 question: What is the cause of	1060
death of Juliet?	1061
Composed question: What is the cause of	1062
death of the person who said a rose	1063
by any other name would smell just	1064
as sweet?	1065
	1066
Hopl question: Who hosted The Price Is	1067
Right before Bob Barker?	1068
Hopl answer: Bill Cullen	1069
Hop2 question: What is the medical	1070
condition of Bill Cullen?	1071
Composed question: What is the medical	1072
condition of the person who hosted	1073
The Price Is Right before Bob	1074
Barker?	1075
2	1010

Hopl question: Who wrote If You Go Away
on a Summer's Day?
Hop1 answer: Rod McKuen
Hop2 question: Which record company
does Rod McKuen own?
Composed question: Which record company
does the person who wrote If You Go
Away on a Summer's Day own?

# A.3 Benchmark Prompt

To use the DETLLM method and generate the final answer, three steps are followed: (1) First-hop prompting, (2) Second-hop prompting, and (3) Final answer generation. The prompt for each stage is provided below. For simplicity, we denote the kretrieved passages as "*Context:* [[1] ... [k]]".

#### Prompt 3: First Hop

Write a search query, query entity, and
complex question
Follow the following format
Context: \${sources that may contain
relevant content}
Question: \${ the question to be answered}
Rationale: Let's think step by step
Based on the context we have
learned the following
\${information from the context that
provides useful clues}
Search Ouery: \${a simple question for
seeking the missing information }
Query Entity: \${query entity name from
search query }
SPARQL: \${SPARQL query used to query
against Wikidata}
Example 1
Context:
Question: What are the occupations of
the person who holds the most
women's Wimbledon titles?
Rationale: Let's think step by step.
Based on the context, we have
learned the following. Decompose
the question to answer the
following single – nop questions. 1.
Wimbledon titles 2 2 What are the
windledon titles? 2. what are the
Search Query: Who holds the most
women's Wimbledon titles?
Ouery Entity: women's Wimbledon titles
SPAROL: None
Example 2
Context:
Question: Which bay is the name of
David Resnick's place of birth?
Rationale: Let's think step by step.
Based on the context, we have
learned the following. Decompose
the question to answer the
following single-hop questions. 1.

Where was David Resnick born? 2. Which bay is the name of this place Search Query: Where was David Resnick born? Query Entity: David Resnick SPARQL: SELECT ?place WHERE {wd:Q962183 wdt:P19 ?place.} Example 3 Context: Question: Is the person who directed the film The Shape of Water a member of the Writers Guild of America, West? Rationale: Let's think step by step. Based on the context, we have learned the following. Decompose the question to answer the following single-hop questions. 1. Who directed the film the shape of water? 2. Is the person the person a member of the Writers Guild of America, West? Search Query: The director of the film The Shape of Water Query Entity: The Shape of Water SPARQL: SELECT ?name WHERE {wd: Q26698156 wdt: P57 ?name.} Target Question Context: Question: How many organizations is the 26th president of the United States a member of? Rationale: Let's think step by step. Based on the context, we have learned the following. Decompose the question to answer the following single-hop questions. 1. who is the 26th president of the United States? 2. How many organizations is this person a member of? Search Query: 26th president of the United States Query Entity: None SPARQL: None 

#### Prompt 4: Second Hop

	1100
Write a search query, query entity, and	1186
SPARQL that will help answer a	1187
complex question.	1188
Follow the following format.	1189
Context:\${sources that may contain	1190
relevant content}	1191
Question: \${ the question to be answered }	1192
Rationale: Let's think step by step.	1193
Based on the context, we have	1194
learned the following.	1195
\${information from the context that	1196
provides useful clues }	1197
Search Query: \${a simple question for	1198
seeking the missing information }	1199
Query Entity: \${query entity name from	1200
search query }	1201
SPARQL: \${SPARQL query used to query	1202
against Wikidata}	1203

1274

Example	1	
Context	• [ [ 1 ]	[1]

Example 1

- Context:[[1] ... [K]] Question: What are the occupations of the person who holds the most women's Wimbledon titles?
- Rationale: Let's think step by step. Based on the context, we have learned the following. Wimbledon is a tennis tournament, and tennis player Martina Navratilova holds the most women's Wimbledon titles. The second step is to answer what are the occupations of this person. Search Query: What are the occupations
- of Martina Navratilova?
- Query Entity: Martina Navratilova
- SPARQL: SELECT ?name WHERE {wd:Q54545
- wdt: P106 ?name. }
- Example 2
- Context:[[1] ... [k]]
- Question: Which bay is the name of David Resnick's place of birth?
- Rationale: Let's think step by step. Based on the context, we have learned the following. David Resnick was born in Rio de Janeiro. The second step is to answer which bay is the name of Rio de Janeiro? Search Query: which bay is the name of Rio de Janeiro?
- Query Entity: Rio de Janeiro SPARQL: None
- Example 3
- Context:[[1] ... [k]]
- Question: Is the person who directed the film The Shape of Water a member of the Writers Guild of America, West?
- Rationale: Let's think step by step. Based on the context, we have learned the following. The Shape of Water is directed by Guillermo del Toro. The second step is to answer is the person a member of the Writers Guild of America, West
- Search Query: the organization Guillermo del Toro is in
- Query Entity: Guillermo del Toro
- SPARQL: SELECT ?name WHERE {wd: Q219124 wdt: P463 ?name. }

Target Question

- Context:[[1] ... [k]]
- Question: How many organizations is the 26th president of the United States a member of?
- Rationale: Let's think step by step. Based on the context, we have learned the following. The 26th president of the United States is Theodore Roosevelt. The second step is to answer how many organizations he is a member of.
- Search Query: How many organizations is Theodore Roosevelt a member of? Query Entity: Theodore Roosevelt SPARQL : SELECT (COUNT(?organization) as ?count) WHERE { wd:Q33866

#### wdt:P463 ?organization. }

#### Prompt 5: Final QA Step

	1270
Answer questions with short factoid	1277
answers.	1278
Follow the following format.	1279
Context: \${ sources that may contain	1280
relevant content}	128
Ouestion: \${ the question to be answered }	128
Rationale: Let's think step by step.	1283
a step -by-step deduction that	1284
identifies the correct response.	128
which will be provided below}	1280
Answer: \${a short factoid answer often	128
hetween 1 and 5 words }	128
between 1 and 5 words j	1280
Example 1	1290
Context: [[1] [k]]	120
Question: What are the occupations of	129
the person who holds the most	1201
women's Wimbledon titles?	120
Rationale: Let's think step by step	120
Martina Navratilova is a tennis	120
player writer poyelist and	120
autobiographer	120
Answer: tennis player writer	1200
novelist and autobiographer	120
noverist, and autobiographer	130
Example 2	130
Context: [[1] [k]]	130/
Ouestion: Which hav is the name of	130
David Respick's place of hirth?	130
Rationale: Let's think step by step	130
David Pasnick was born in Pio de	1300
Laneiro and "Rio de Laneiro" was	130
the name of Guanahara Bay	1300
Answer: Guanabara Bay	131(
Answer. Guanabara Day	131
Example 3	131
Context · [[1] [k]]	1313
Question: Is the person who directed	131/
the film The Shape of Water a	131
member of the Writers Guild of	1310
America West?	1313
Rationale: Let's think step by step.	1318
Guillermo del Toro Gomez is a	1319
filmmaker, he is a member of the	1320
Writers Guild of America, West	132
Answer: yes	132
	1323
Target	1324
Context:[[1] [k]]	132
Ouestion: How many organizations is the	1320
26th president of the United States	132
a member of?	1328
Rationale: The 26th president of the	1329
United States was Theodore	1330
Roosevelt. He is a member of 5	133
organizations.	133
Answer: 5	133
	1000

### A.4 Human Annotation Instruction

We show the instructions and annotating exam-1335 ples provided to human annotators to annotate the dataset as below. 1337

1275

Bridge Entity: LeBron James is to judge and revise the complex question chained 1339 Composed Question: How many children by two single-hop questions. To complete this goal, 1340 does the highest-paid athlete in you need to do the following two tasks: the NBA have? 1341 An accepted question should meet the following • Judge and revise the single-hop question gen-1342 criteria: erated from the knowledge base triplet. 1343 The composed question must be constructed Judge and revise the composed complex ques-1344 1345 tion. using two single-hop questions, with the answer to the first question becoming the subject 1346 **Task 1** Given a triplet (subject, relation, object) of the second question. and a machine-generated question as shown below, 1347 you need to judge the quality of the generated ques-1348 • Ensure that the composed question does not tion and whether it is acceptable, needs revision, or 1349 reveal the answer itself. is rejected. If the question can be revised, please re-1350 • Use 'Answer 2' as the answer to the composed 1351 vise the question rather than reject it. If the question question. is too poor to revise, reject the question. 1352 1353 Triplet: (LeBron James; child; [Bryce If you choose to reject the question, please select James, Zhuri James, Bronny James]) 1354 one of the following reasons. If your reason is not Question: How many children does LeBron 1355 1356 James have? listed, choose 'Other' and include a comment. 1357 An accepted triplet question should satisfy the • Circular question: Two single-hop questions following criteria: are the same question. • The question focuses on the subject w.r.t rela-1359 Bridge entity answer leaking. tion. 1360 Final answer leaking. • The question should sound natural and fluent. • The answer to the generated question should 1362 • Change in the original meaning of single-hop be the object, thus the object cannot be shown 1363 questions. in the question. 1364 • Other. **Task 2** Judge and revise the composed complex 1365 question given the following information. If the A.5 An overview of DETLLM data 1366 question can be revised, please revise the question generation process 1367 rather than reject it. If the question is too poor to 1368 An overview of DETLLM data generation process revise, reject the question and choose the reason is shown in Figure 3. for rejection. 1370 Below is a list of provided information: 1371 A.6 Dataset Analysis The stats of the dataset are shown in Figure 4. • Two single-hop question-answer pairs: 1372 "(Question 1, Answer 1)" and "(Question 2, 1373 A.7 **Anecdotal Examples for Representative** Answer 2)". Types • The bridging entity "Bridge Entity" that 1375 Anecdotal Examples for Representative Types are chains two single-hop questions together. 1376 shown in Table 5. Machine generated composed question "Com-1377 A.8 Additional Experiment Results 1378 posed Question". SPARQL Generation Analysis Symbolic language generation is an essential tool, which is ex-Question 1: Who is the highest-paid 1379 1380 athlete in the NBA ecuted against the Wikidata engine to assist with 1381 Answer 1: LeBron James structured knowledge retrieval. We provide a de-1382 Question 2: How many children does tailed breakdown analysis of SPARQL generation 1383 LeBron James have?

Answer 2: [3, three]

1384

1385

1386

1387

1388

1389

1391

1392

1393

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

**Overall Instruction** The goal of the annotation



a. Types of questions.

Figure 4: Types of (a) questions, and (b) KB relations, covered in DIVKNOWQA.

in Table 6. "QID" represents the percentage of ex-1425 1426 amples with entity IDs correctly linked to Wikidata. Additionally, we present the percentage of 1427 examples linked to the Wikidata in terms of both 1428 entity IDs and relation IDs denoted as "QID+REL". 1429 The last column, labeled "QID\*", showcases the 1430 percentage of examples with great potential for ac-1431 curate identification through entity disambiguation. 1432 In our experimental process, we first identify the 1433 entity name from the decomposed question as a re-1434 triever query and then link the entity from the query 1435 1436 to Wikidata. The returned results provide a list of candidate Wikidata entities, from which we select 1437 the most semantically similar one by computing 1438 the similarity between the query and the entity's 1439 description. The displayed number reveals that this 1440 heuristic entity disambiguation process fails to rec-1441

ognize those examples that actually contain the<br/>correct entity ID within the candidate list. This<br/>highlights a potential avenue for further improving<br/>model performance.1442<br/>1443

**Establishing Oracle Performance** In Table 7, 1446 we present the experimental results obtained using 1447 Oracle information. In these experiments, we grant 1448 the model access to ground-truth passages from the 1449 Oracle Text and linearized KB triplets from the KB 1450 Oracle. A notable observation is the comparison 1451 between Text Oracle and KB Oracle. We find that 1452 KB Oracle exerts a more significant influence on 1453 the final results. This is because structured knowl-1454 edge contains long-tail knowledge, showing the 1455 necessity to effectively explore structured knowl-1456 edge. Furthermore, when both Text and KB Ora-1457 cle sources are provided, the model's performance 1458 Table 5: Types of multi-hop reasoning required to answer questions in DIVKNOWQA. Two single-hop questions are shown: TextQA is sampled from NQ, and KBQA is generated using the sampled KB-Triplet. The question from **DIVKNOWQA** is based on those two single-hop questions.

Order	Туре	%	Example				
Text → KB	short entity	20.3	TextQA: Who is Rafael Nadal married to?         Answer: María Francisca Perelló         KB-Triplet: (Rafael Nadal, spouse, María Francisca Perelló)         KBQA: Who won the Men's US Open 2017?         Answer: Rafael Nadal         DIVKNOWQA: Who is the person married to the winner of the Men's US Open 2017?         Answer: María Francisca Perelló				
	yes/no	17.9	TextQA: Who sang When the Lights Went Out in Georgia?         Answer: Vicki Lawrence         KB-Triplet: (Vicki Lawrence, hair color, red hair)         KBQA: What is Vicki Lawrence's hair color?         Answer: red hair         DIVKNOWQA: Is the hair color of the singer of "When the Lights Went Out in Georgia" gray?         Answer: no				
	aggregate	21.1	TextQA: who does Meg 's voice on Family Guy?         Answer: Vicki Lawrence         KB-Triplet: (Mila Kunis, child, [Wyatt Kutcher, Dimitri Kutcher])         KBQA: How many children does Mila Kunis have?         Answer: Two         DIVKNOWQA: How many children does the person who does Meg's voice on Family Guy have?         Answer: Two				
$\mathbf{KB}  ightarrow \mathbf{Text}$	short entity	20.7	KB-Triplet: (William Weatherall Wilkins, present in work, Mary Poppins Returns)         KBQA: In which work is William Weatherall Wilkins present?         Answer: Mary Poppins Returns         TextQA: Who play Mary Poppins in Mary Poppins Returns?         Answer: Emily Blunt         DIVKNOWQA: Who plays Mary Poppins in the work in which William Weatherall Wilkins is present?				
	yes/no	20.0	KB-Triplet: (Girl #2, present in work, High School Musical)         KBQA: In which work is Girl #2 present?         Answer: High School Musical         TextQA: What grade were they in in high school musical 1?         Answer: juniors         DIVKNOWQA: Were they seniors in the work in which Girl #2 is present?         Answer: no				

Table 6: Breakdown Analysis on SPARQL generation.

	QID	QID+REL	QID*
Text-KB(Sparse)	26.5	22.4	6.91
Text-KB(Dense)	31.8	27.6	7.87
Text(Sparse)-KB(Sparse)	26.3	22.6	7.34
Text(Dense)-KB(Sparse)	29.7	26.5	7.66

Table 7: Experiment results using Oracle knowledgesource retrieval in each sub-step.

	EM	F1	Recall	H1-R	H2-R
Oracle_Text	26.8	31.3	33.4	96.4	51.7
Oracle_KB	38.1	40.0	42.2	100.0	62.1
Oracle_All	48.7	52.2	52.8	100.0	96.7

reaches an Exact Match (EM) rate of 48.7%, highlighting the necessity of each knowledge source. In comparison to our current established results from DETLLM, this benchmark reveals substantial room for the research community to further explore and improve upon.