On the Robustness of RAG Systems in Educational Question Answering under Knowledge Discrepancies

Anonymous ACL submission

Abstract

Retrieval-Augmented Generation (RAG) systems have demonstrated remarkable potential as question answering systems in the K-12 Education domain, where knowledge is typically queried within the restricted scope of authoritative textbooks. However, the discrepancy between textbooks and the parametric knowledge in Large Language Models (LLMs) could undermine the effectiveness of RAG systems. To systematically investigate the robustness of RAG systems under such knowledge discrepancies, we present EDUKDQA, a question answering dataset that simulates knowledge discrepancies in real applications by applying hypothetical knowledge updates in answers and source documents. EDUKDQA includes 3,005 questions covering five subjects, under a comprehensive question typology from the perspective of context utilization and knowledge integration. We conducted extensive experiments on retrieval and question answering performance. We find that most RAG systems suffer from a substantial performance drop in question answering with knowledge discrepancies, while questions that require integration of contextual knowledge and parametric knowledge pose a challenge to LLMs. All resources will be released to foster further research.

1 Introduction

017

023

024

040

043

In K-12 education, Question Answering (QA) systems serve as an important resource for learning assistance, where answers are precisely provided within a restricted knowledge scope from authoritative sources (i.e., textbooks) (Raamadhurai et al., 2019; Soares et al., 2021). Meanwhile, benefiting from the emergent abilities (Wei et al., 2022) of Large Language Models (LLMs) and advanced information retrieval (IR) methods, Retrieval-Augmented Generation (RAG) systems have achieved remarkable performance in various knowledge-intensive tasks in natural language processing (Lewis et al., 2020; Jiang et al., 2023; Gao



Figure 1: An illustration of knowledge discrepancy in educational QA and the application of RAG systems.

et al., 2024), demonstrating their great potential as QA systems in K-12 education (Gan et al., 2023; Kasneci et al., 2023; Yan et al., 2024).

In K-12 educational QA, one of the primary concerns is ensuring that the knowledge conveyed in the answer is consistent with the officially designated textbooks (Extance, 2023). However, there are notable discrepancies between the knowledge in textbooks and the internal knowledge of LLMs, due to the evolving nature of facts (Arbesman, 2012), updates in pedagogical approaches (Provenzo et al., 2011), as well as regional and cultural variations in content (Patel, 2015). It remains unclear whether RAG systems can robustly incorporate knowledge from authoritative sources and generate consistent answers under such knowledge discrepancies (sce-

Question Type	Reasoning Pattern and Example Question	Hypothetical Knowledge Update
Simple Direct	○→◇→●	NV goggles - detect - Infrared light <i>Ultraviolet</i>
Simple Direct	What type of light is <i>detected</i> by night vision goggles ?	Original: Infrared Light Updated: Ultraviolet Light
Multi-hon Direct	$\bigcirc \rightarrow \diamond \rightarrow \bigcirc \rightarrow \diamond \rightarrow \bigcirc$	de Broglie eq developed by - Louis de Broglie <u>Maurice</u>
тап пор Блеет	Which scientist <i>developed</i> an equation that can <i>calculate</i> the wavelength of a particle ?	Original: Louis de Broglie Updated: Maurice de Broglie
Multi-hop Distant		$\frac{Na^+/K^+}{Ca^{2+}}$ Pump - creates - EC Gradient
·	Which pump <i>creates</i> an electrochemical gradient that <i>enables</i> secondary active transport to occur?	Original: Sodium-potassium Pump Updated: Calcium Pump
Multi hon Implicit	○→◇→○ ···◇···◇	Polonium - found in - Uranium ores <i>Thorium</i>
миш-пор ітриси	Who <i>discovered</i> the radioactive element that is commonly <i>found</i> in uranium ores ?	Original: Marie Curie Updated: Jöns Jacob Berzelius
Distant Implicit		Mitochondrion - conducts - Cellular respiration Golgi apparatus
	Who <i>discovered</i> the organelle that is <i>responsible</i> for the bio-logical process that <i>produces</i> ATP ?	Original: Albert von Kölliker Updated: Camillo Golgi
O su	bject/object 🔷 predicate 🔘 🔷 contextual fact 🔵 updated fact	parametric fact <i>dist.</i> distant fact

Table 1: Five question types in EDUKDQA with their reasoning patterns, example questions, and hypothetical knowledge update illustrated in factual triplets and answers. In example questions, subjects and objects are marked in **bold**, while predicates are marked in *italic*.

nario illustrated in Figure 1).

To fill this gap, we aim to systematically assess the robustness of RAG systems in performing question answering in K-12 education when encountering knowledge discrepancies. We present EDUKDQA (Educational Knowledge Discrepancy Question Answering), a new dataset containing 3,005 multiple-choice questions covering the subjects of Physics, Chemistry, Biology, Geography, and History from the middle-school curriculum. To simulate the knowledge discrepancy between LLMs and textbooks, we conduct a hypothetical knowledge update, in which we modify the original factual knowledge from the textbooks into plausible alternatives while maintaining coherent and consistent context. Moreover, we tailored a comprehensive question typology to stress-test the context utilization and knowledge integration abilities of LLMs under such scenarios.

We conducted extensive experiments with various retrieval methods and LLMs. We find that most RAG systems still suffer from a considerable performance drop when facing knowledge discrepancies. Notably, while most LLMs can incorporate distant contextual facts well, they struggle to integrate their parametric knowledge with contextual knowledge effectively. In terms of retrieval, traditional lexical-based methods show advantages due to specificity of the academic terms, and their performance may be further enhanced through an ensemble reranking mechanism. To encourage further research, we will make our benchmark dataset and the associated code publicly available.

090

091

095

096

097

2 The EDUKDQA Dataset

In this section, we introduce the methodology of hypothetical knowledge update and the design of our question typology. The curation pipeline and detailed statistics of the EDUKDQA dataset are provided in Appendix B and E, respectively.

2.1 Hypothetical Knowledge Update

One of the core objectives of our dataset is to simulate the knowledge discrepancy between LLMs 101 and authoritative textbooks when performing edu-102 cational QA. However, such discrepancies are often 103 fuzzy and highly sparse in real-world data, mak-104 ing it infeasible to collect and organize. Consequently, we designed the methodology of hypo-106 thetical knowledge update, performing it on high-107 quality open-source textbooks. The general proce-108 dures are as follows: 1) Curate factual questions 109 from textbook paragraphs following our designed 110 question typology. 2) Select a plausible but factu-111 ally incorrect answer as the updated ground-truth 112 answer. 3) Replace all occurrences of the origi-113 nal answer in the paragraph with the updated an-114 swer and adjusted other relevant statements in the 115

context to ensure that the updated paragraph is
coherent and consistent. This process is further
guaranteed through extensive human curation and
verification. Examples of hypothetical knowledge
updates in our dataset are provided in Appendix D.

2.2 Question Typology

121

122

123

124

125

126

127

128

129

130

131

132

133

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

152

153

154

156

157

158

160

161

We identify two potential challenges for LLMs when performing QA with knowledge discrepancy: 1) Context Utilization: Whether LLMs can identify and utilize the corresponding facts from the context; and 2) Knowledge Integration: Whether LLMs can incorporate their own parametric knowledge with the contextual ones in question answering. To this end, we designed our question typology, as illustrated in Table 1, to investigate such abilities. Based on the two basic question types, Simple Direct and Multi-hop Direct, we developed the Multi-hop Distant type to evaluate the context utilization ability for distant facts from the passage, and the *Multi-hop Implicit*¹ type to evaluate the knowledge integration ability that combines their own factual knowledge with retrieved ones. Moreover, the Distant Implicit type poses a greater challenge by combining both features. To ensure our evaluation of knowledge integration ability is independent of knowledge coverage, we restrict the facts requiring LLMs' own knowledge to be high-frequency only (Sun et al., 2024).

3 Experiments and Analyses

Typically, RAG systems first conduct document retrieval based on given queries, then perform question answering with LLMs based on the retrieved information loaded in the context. In this section, we comprehensively evaluate and analyze the performance of retrieval methods and LLMs on the EDUKDQA dataset. For details of all tested methods and models, please refer to Appendix C.

3.1 Retrieval Performance

Experimental results of the retrieval methods are presented in Table 2. Traditional lexical retrieval methods, such as BM25, demonstrated strong performance on our dataset, while dense retrieval methods, such as Mistral-embed and Ada-002, achieved comparable performance. Since our dataset focuses on the K-12 education domain, lexical retrieval effectively captures domain-specific

Retrieval Methods	Category	R@1	R@5
TF-IDF (Spärck Jones, 1972)	Lexical	65.82	88.72
BM25 (Robertson et al., 1994)	Lexical	82.73	95.27
SPLADE (Formal et al., 2021)	Lexical/Dense	78.04	90.12
Contriever (Izacard et al., 2022)	Dense	53.18	81.80
Conmsmarco (Izacard et al., 2022)	Dense	76.17	93.54
Mistral-embed (Mistral, 2023a)	Dense	78.74	95.31
Ada-002 (OpenAI, 2022)	Dense	79.23	95.44
Query Rewrite (Ma et al., 2023)	Pre-Retrieval	78.87	94.21
Hybrid Rerank (BM25 + Ada-002)	Ensemble	84.43	96.04
Contriever (fine-tuned)	Dense+FT	84.19	98.96
Conmsmarco (fine-tuned)	Dense+FT	87.95	99.50

Table 2: Document retrieval performances (in recall @ 1/5 %) of retrieval methods from different categories in the EDUKDQA dataset. The highest scores are marked as **bold**, while the 2nd and 3rd-best scores are <u>underlined</u>. The retrieval granularity is set to paragraph.

keywords in the queries and identifies the corresponding documents. This characteristic underscores the need for fine-tuning dense retrieval models on our dataset. To this end, we fine-tuned both Contriever and Contriever-msmarco on our documents, resulting in significant improvements and highlighting the importance of corpus-specific finetuning in educational document retrieval. 162

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

186

187

188

189

191

192

193

195

Moreover, ensemble methods have demonstrated their effectiveness in various information retrieval tasks (Thakur et al., 2021). We implemented a hybrid approach that retrieves the top k^2 documents with Ada-002, followed by re-ranking with BM25. This ensemble method marginally improved the retrieval performance in both metrics. Conversely, query rewriting³ (Ma et al., 2023) implemented upon BM25 did not yield performance gains. Detailed results across subjects and question types are provided in Appendix F.

3.2 Question Answering Performance

Experimental results of LLMs in question answering are presented in Table 3. Most LLMs exhibited comparable performance for question types *Simple Direct, Multi-hop Direct,* and *Multi-hop Distant.* This suggests that **both multi-hop reasoning and distant context utilization do not pose** *significant challenges for modern LLMs.* However, for *Multi-hop Implicit* questions requiring the integration of contextual and internal knowledge, a substantial performance disparity emerges between smaller open-source LLMs and advanced models. This disparity is further amplified for *Distant Implicit* questions. The most capable model, o1-preview, attains over 80% accuracy on *Implicit* questions, whereas Mistral-7b's performance

¹The word 'implicit' indicates that the questions indirectly query the updated facts by embedding them within the middle of the multi-hop reasoning chain.

²We use k = 6 as the optimal hyperparameter setting.

³We use Mistral-small-2409 as the LLM query rewriter.

Large Language Models	Question Typology									
	Simple Direct	Multi-hop Direct	Multi-hop Distant	Multi-hop Implicit	Distant Implicit					
Mistral-7b (Mistral, 2023b)	77.70	69.31	72.74	45.32	33.98	61.26				
Mixtral-8x22b (Mistral, 2023c)	84.10	84.58	87.15	73.86	65.37	79.67				
Mistral-small-2409 (Mistral, 2024a)	87.87	88.70	89.69	72.18	60.81	80.77				
Mistral-large-2407 (Mistral, 2024b)	83.93	83.51	87.29	82.25	70.57	81.66				
Gemini-1.5-flash (Google, 2024a)	80.98	82.44	87.57	76.02	63.25	78.54				
Gemini-1.5-pro (Google, 2024b)	86.56	87.63	88.42	76.26	63.58	81.10				
Llama3-8b (Meta, 2024)	90.33	88.85	89.69	63.55	49.43	77.77				
Llama3-70b (Meta, 2024)	<u>96.72</u>	96.49	96.89	79.14	63.25	87.42				
GPT-3.5-turbo (OpenAI, 2022)	92.62	90.53	91.24	71.22	57.72	81.73				
GPT-4 (OpenAI, 2023a)	89.51	89.47	90.68	78.66	70.89	84.46				
GPT-4-turbo (OpenAI, 2023b)	95.74	96.18	96.19	81.06	71.71	88.99				
GPT-40 (OpenAI, 2024a)	91.97	94.81	93.64	81.29	70.41	87.09				
Claude-3-sonnet (Anthropic, 2024a)	94.59	92.82	93.64	77.70	60.65	84.69				
Claude-3.5-sonnet (Anthropic, 2024b)	97.54	96.49	95.62	83.69	73.82	90.08				
o1-mini (OpenAI, 2024b)	95.90	95.73	97.03	85.85	75.45	90.55				
o1-preview (OpenAI, 2024b)	95.08	97.71	97.46	86.33	78.86	91.68				

Table 3: Question answering performances (in accuracy %) of LLMs in the EDUKDQA benchmark with corresponding documents provided. The highest scores are marked as **bold**, while the 2^{nd} and 3^{rd} -best scores are <u>underlined</u>. All LLMs are tested under zero-shot settings, with a *Locate-and-Answer* prompting approach that facilitates active information acquisition from contextual documents, with details in Appendix G.

RAG System	Hypothetical	Drop	
·	Before	After	•
Llama3-8b + Ada-002 Llama3-8b + Rerank GPT-4o + Ada-002 GPT-4o + Rerank	87.49 88.49 96.57 97.10	62.60 66.02 69.65 73.71	24.89 22.47 26.92 23.39

Table 4: Performance drop of RAG systems with hypothetical knowledge updates in our benchmark.

falls below 40%. These findings indicate that
knowledge integration is an emergent capability
presenting greater difficulties for LLMs under
knowledge discrepancies, particularly when coupled with the need for distant context utilization.

3.3 Overall Performance

How do knowledge discrepancies affect the performance of RAG systems in educational application?
We selected two representative LLMs: Llama3-8b from open-source models and GPT-4o for proprietary models, combined with two high-performing retrieval methods: Ada-002 and hybrid rerank. The resulting RAG systems were tested on questions *before* and *after* hypothetical knowledge updates, with results presented in Table 4. We observed an accuracy drop of 22-27%, indicating a substantial performance degradation in modern RAG systems when faced with knowledge discrepancies.

4 Related Work

217Retrieval-Augmented GenerationFollowing218the categorization by Gao et al. (2024), the219RAG methods in our experiment are from Naive220RAG (lexical, dense) and Advanced RAG (rerank,

rewrite). Recently, *Modular RAG* has emerged to enhance the adaptability and versatility of RAG systems (Shao et al., 2023; Asai et al., 2023).

Educational Question Answering Prior to the emergence of RAG systems utilizing LLMs, various educational QA systems were developed to provide pedagogically appropriate responses to student inquiries (Abdi et al., 2018; Agarwal et al., 2019). While recent literature explores LLM applications in QA and learning assistance roles (Nye et al., 2023; Kuo et al., 2023; Wang et al., 2024).

Knowledge Discrepancy in LLMs Mitigating knowledge discrepancies or conflicts in applications is a fundamental challenge in LLM research (Xu et al., 2024). Researchers have proposed tuning-based (Li et al., 2022) and prompting-based (Zhou et al., 2023) methods to enhance LLMs' robustness under such conflicts. However, these conflicts in educational applications and RAG systems remain relatively underexplored.

5 Conclusion

This paper systematically evaluates the robustness of RAG systems in K-12 educational question answering under knowledge discrepancies using a comprehensive dataset. Experimental findings reveal substantial performance degradation in RAG systems when faced with knowledge discrepancies, which is primarily attributed to deficiencies in incorporating contextual and parametric knowledge in question answering—an emergent and challenging ability for modern large language models.

Limitations

257

261

262

264

271

281

287

290

We discuss three main limitations of our work.

First, EDUKDQA employs the approach of hypothetical knowledge updates to effectively simulate real-world knowledge discrepancies for two primary reasons: (1) Real-world knowledge conflicts are often sparse, noisy, and difficult to systematically collect or organize into a cohesive dataset 259 suitable for large-scale evaluation. Hypothetical 260 updates provide a scalable and controlled alternative, allowing us to bypass these limitations. (2) By leveraging a systematic annotation and curation 263 pipeline, we can generate questions with diverse and well-defined reasoning patterns that align with 265 our typology, enabling more robust evaluation of complex question-answering tasks. For future research centered on real-world knowledge discrepancies, we recommend an alternative methodology 269 that incorporates temporal attributes (Chen et al., 270 2021; Zhang and Choi, 2024). This approach focuses on identifying outdated facts and capturing time-sensitive data (e.g., economic trends, annual events, or societal changes) to construct datasets that reflect real-world knowledge updates. While promising, this method is constrained by the lim-276 ited overlap between time-sensitive data-often numerical or attribute-specific-and the broader contextual needs of educational question-answering tasks, which may reduce the comprehensiveness of the resulting datasets.

> Next, regarding document retrieval, some recent hierarchical retrieval paradigms, such as GraphRAG (Edge et al., 2024) and HippoRAG (Gutiérrez et al., 2024), are not included in our experiments due to their implementation complexity. However, we believe that such structured paradigms could effectively enhance retrieval performance in our scenario, as educational documents are well-structured and contain high-quality factual knowledge.

Finally, this paper primarily evaluates the robustness of RAG systems in the proposed scenario, with experiments conducted using various retrieval methods and large language models. Potential improvements could be achieved through the design of tailored reasoning frameworks via prompting, in-context learning or alignment in LLMs, which we leave for future research.

Ethics Statement

In constructing the EDUKDQA dataset, we col-301 lected text from open-access textbooks. Detailed 302 sources and licenses are provided in Appendix A. 303 The human curation and verification in our annota-304 tion pipeline were carried out by a group of post-305 graduate students with extensive experience in NLP 306 research. We ensured that the updated knowledge 307 is free from harmful or toxic content. It is impor-308 tant to note that our dataset is designed solely to 309 evaluate the robustness of Retrieval-Augmented 310 Generation systems under scenarios with knowl-311 edge discrepancies and is not suitable for assessing 312 the factual accuracy of QA systems. 313

300

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

References

- Asad Abdi, Norisma Idris, and Zahrah Ahmad. 2018. Qapd: an ontology-based question answering system in the physics domain. Soft Computing, 22:213-230.
- A. Agarwal, N. Sachdeva, R. K. Yadav, V. Udandarao, V. Mittal, A. Gupta, and A. Mathur. 2019. Eduqa: Educational domain question answering system using conceptual network mapping. In ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 8137-8141.
- Anthropic. 2024a. Blog: Introducing the next generation of claude.
- Anthropic. 2024b. Introducing claude 3.5 sonnet.
- S. Arbesman. 2012. The Half-life of Facts: Why Everything We Know Has an Expiration Date. Human body. Current.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection.
- Wenhu Chen, Xinyi Wang, and William Yang Wang. 2021. A dataset for answering time-sensitive questions.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization.
- Andy Extance. 2023. ChatGPT has entered the classroom: how LLMs could transform education. Nature, 623:474-477.
- Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. Splade v2: Sparse lexical and expansion model for information retrieval.

- Wensheng Gan, Zhenlian Qi, Jiayang Wu, and Jerry Chun-Wei Lin. 2023. Large language models in education: Vision and opportunities.
 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,
 - Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-augmented generation for large language models: A survey.
 - Google. 2024a. Gemini flash google deepmind.
- Google. 2024b. Gemini pro google deepmind.

356

360

366

367

369

373

375 376

379

386

395

398

400

401

402

- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2024. Hipporag: Neurobiologically inspired long-term memory for large language models.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning.
- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation.
- Enkelejda Kasneci, Kathrin Sessler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, Stephan Krusche, Gitta Kutyniok, Tilman Michaeli, Claudia Nerdel, Jürgen Pfeffer, Oleksandra Poquet, Michael Sailer, Albrecht Schmidt, Tina Seidel, Matthias Stadler, Jochen Weller, Jochen Kuhn, and Gjergji Kasneci. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and Individual Differences*, 103:102274.
- Bor-Chen Kuo, Frederic T. Y. Chang, and Zong-En Bai. 2023. Leveraging llms for adaptive testing and learning in taiwan adaptive learning platform (talp). In *Empowering Education with LLMs – the Next-Gen Interface and Content Generation*, volume 3487, pages 1–14, Tokyo, Japan. CEUR Workshop Proceedings, CEUR-WS.org.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, and Sanjiv Kumar. 2022. Large language models with controllable working memory. *arXiv preprint arXiv:2211.05110*.
- Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting for retrievalaugmented large language models.
- Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date.

Mistral. 2023a. Embeddings mistral ai large language models.	403 404
Mistral. 2023b. Mistral 7b frontier ai in your hands.	405
Mistral. 2023c. Mixtral of experts frontier ai in your hands.	406 407
Mistral. 2024a. Ai in abundance mistral ai frontier ai in your hands.	408 409
Mistral. 2024b. Large enough mistral ai frontier ai in your hands.	410 411
Benjamin D. Nye, Dillon Mee, and Mark G. Core. 2023. Generative large language models for dialog-based tutoring: An early consideration of opportunities and concerns. In <i>Empowering Education with LLMs – the</i> <i>Next-Gen Interface and Content Generation</i> , volume 3487, pages 1–15, Tokyo, Japan. CEUR Workshop Proceedings, CEUR-WS.org.	412 413 414 415 416 417 418
OpenAI. 2022. Chatgpt: Optimizing language models for dialogue.	419 420
OpenAI. 2022. New and improved embedding model.	421
OpenAI. 2023a. Gpt-4 technical report.	422
OpenAI. 2023b. New models and developer products announced at devday.	423 424
OpenAI. 2024a. Hello gpt-4o.	425
OpenAI. 2024b. Introducing openai o1.	426
Leigh Patel. 2015. <i>Decolonizing Educational Research:</i> <i>From Ownership to Answerability</i> , 1 edition. Rout- ledge, New York.	427 428 429
E.F. Provenzo, A.N. Shaver, and M. Bello. 2011. <i>The</i> <i>Textbook as Discourse: Sociocultural Dimensions of</i> <i>American Schoolbooks</i> . Taylor & Francis.	430 431 432
Srikrishna Raamadhurai, Ryan Baker, and Vikraman Poduval. 2019. Curio SmartChat : A system for nat- ural language question answering for self-paced k-12 learning. In <i>Proceedings of the Fourteenth Workshop</i> <i>on Innovative Use of NLP for Building Educational</i> <i>Applications</i> , pages 336–342, Florence, Italy. Associ- ation for Computational Linguistics.	433 434 435 436 437 438 439
Stephen E Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. 1994. Okapi at trec-3. In Overview of the Third Text REtrieval Conference (TREC-3). NIST.	440 441 442 443
Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. En- hancing retrieval-augmented large language models with iterative retrieval-generation synergy.	444 445 446 447
Teotino Gomes Soares, Azhari Azhari, Nur Rokhman, and E Wonarko. 2021. Education question answering systems: a survey. In <i>Proceedings of The Interna</i> -	448 449 450

tional MultiConference of Engineers and Computer

451

452

6

Scientists.

Karen Spärck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.

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

482

483

484

485

486

487

488

489

490

491

492

493

- Kai Sun, Yifan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. 2024. Head-to-tail: How knowledgeable are large language models (LLMs)? A.K.A. will LLMs replace knowledge graphs? In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 311–325, Mexico City, Mexico. Association for Computational Linguistics.
- Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. Beir:
 A heterogenous benchmark for zero-shot evaluation of information retrieval models.
- Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S. Yu, and Qingsong Wen. 2024. Large language models for education: A survey and outlook.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models.
- Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024. Knowledge conflicts for llms: A survey.
- Lixiang Yan, Lele Sha, Linxuan Zhao, Yuheng Li, Roberto Martinez-Maldonado, Guanliang Chen, Xinyu Li, Yueqiao Jin, and Dragan Gašević. 2024. Practical and ethical challenges of large language models in education: A systematic scoping review. *British Journal of Educational Technology*, 55(1):90– 112.
- Michael J. Q. Zhang and Eunsol Choi. 2024. Mitigating temporal misalignment by discarding outdated facts.
- Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2023. Context-faithful prompting for large language models.



Figure 2: Distribution of five subjects and their corresponding topics included in the EDUKDQA dataset.

A Source Textbooks

The documents in the EDUKDQA dataset is organized based on the following public textbooks: 494

495

496

497

498

499

501

502

503

504

505

506

508

509

510

511

512

513

Physics *Physics* by Openstax (CC-BY-4.0⁴) https://openstax.org/details/books/ physics

Chemistry *Chemistry* by Openstax (CC-BY-4.0¹) https://openstax.org/details/books/ chemistry-2e

Biology *Biology* by Openstax (CC-BY-4.0¹) https://openstax.org/details/books/ biology-2e

History *World History* by OER Commons (CC-BY-NC-4.0⁵) https://oercommons.org/ courses/world-history-2

Geography World Regional Geography by Sailor Academy (CC-BY-3.0⁶) https://learn.saylor. org/course/view.php?id=722

Meanwhile, the detailed topics included in each subjects are illustrated in Figure 2.

⁴https://creativecommons.org/licenses/by/4.0/deed.en ⁵https://creativecommons.org/licenses/by-nc/4.0/deed.en

⁶https://creativecommons.org/licenses/by/3.0/



Figure 3: An overview of the data curation pipeline of the EDUKDQA dataset.

B Curation Pipeline

514

515

516

517

518

519

524

526

528

531

535

539

540

541

543

545

549

The curation pipeline of our dataset is illustrated in Figure 3. We first perform triplet extraction on the textbook documents and generate a documentlevel knowledge graph (KG). Next, we perform sub-graph matching based on fixed reasoning patterns to sample candidate queries, and selectively transform them into natural language questions. Then, hypothetical knowledge update is executed and verified to guarantee consistency between the updated answer and the document.

Context-focused question types, including *Simple Direct*, *Multi-hop Direct*, and *Multi-hop Distant*, are acquired through this process. For *Multi-hop Distant* questions, we leverage distant facts, defined as connected triplet pairs that are separated in the document's sequential ordering. These questions are only assigned to documents containing no fewer than 200 words. To generate the other two *Implicit* question types that require knowledge integration, we perform extra QA augmentation followed by an expert verification process.

Our data curation process is performed through an integrated framework involving both human annotators and LLMs. For LLM annotation, we adopted Claude-3.5-Sonnet for its outstanding instruction-following ability. Following manual verification, 90.5% of these queries were retained or underwent minor refinements to become highquality questions, yielding an overall success rate of 86.4%. The total API cost for data annotation is approximately 300 USD.

C Model Details

In this section, we briefly introduce all tested retrieval methods and large language models in our benchmark experiment.

In the retrieval stage of our experiments, we employ a diverse range of retrieval methods. For traditional lexical retrieval, we include TF-IDF (Spärck Jones, 1972), which vectorizes documents and queries based on term frequency and inverse document frequency, and BM25 (Robertson et al., 1994), which enhances TF-IDF with document length normalization and probabilistic term weighting. We also incorporate SPLADE (Formal et al., 2021), a method that bridges lexical and dense retrieval paradigms. For dense retrieval, we evaluate several methods that encode questions and documents into the same vector space: Contriever (Izacard et al., 2022), an unsupervised text encoder; its fine-tuned variant Contriever-msmarco, which we further enhanced by applying contrastive learning (Izacard et al., 2022) to fine-tune both models on the EDUKDQA dataset for improved retrieval capability; and two closed-source embedding models, Mistral-embed (Mistral, 2023a) and Ada-002 (OpenAI, 2022). Additionally, we explore a preretrieval method, Query Rewrite (Ma et al., 2023), which reformulates queries to improve retrieval performance. Finally, we implement a Hybrid Rerank approach that combines BM25 and Ada-002, leveraging the strengths of both lexical and dense retrieval methods.

550

551

552

553

554

555

556

557

558

559

560

561

562

563

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

586

In our evaluation, we employ a diverse set of state-of-the-art large language models to assess their performance across various tasks. The models include: Mistral AI's open-source models, ranging from the compact **Mistral-7b** (Mistral, 2023b) to the more advanced MoE model **Mixtral-8x22b** (Mistral, 2023c), and their latest iterations **Mistral-small-2409** (Mistral, 2024a) and **Mistrallarge-2407** (Mistral, 2024b). Google's Gemini models are represented by **Gemini-1.5-flash**

633 634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

(Google, 2024a) and Gemini-1.5-pro (Google, 587 2024b). Meta's Llama3 series is included with 588 both 8b and 70b parameter versions (Meta, 2024). 589 We also evaluate OpenAI's models, including GPT-3.5-turbo (OpenAI, 2022), GPT-4 (OpenAI, 2023a), GPT-4-turbo (OpenAI, 2023b), and GPT-592 40 (OpenAI, 2024a). Anthropic's LLMs are repre-593 sented by Claude-3-sonnet (Anthropic, 2024a) and Claude-3.5-sonnet (Anthropic, 2024b). Lastly, we include OpenAI's latest o1 series, o1-mini and o1preview (OpenAI, 2024b), which achieve remarkable performance across various metrics through 598 inference-time scaling. This comprehensive selec-599 tion allows us to compare a wide range of model architectures and sizes, providing valuable insights into the current state of LLM capabilities.

D Dataset Example

We provide an example question and its corresponding paragraph before and after hypothetical knowledge update in Table 5.

E Dataset Statistics

The statistics for the documents and questions in the EDUKDQA dataset are provided in Table 6 and Table 7, respectively.

F Result Details

606

607

610

611

612

613

614

615

616

617

619

620

621

The comprehensive evaluation results for retrieval approaches and large language models across various subjects and question types are presented in Table 8 and Table 9, respectively. We observed that the performance of lexical retrieval methods correlates positively with question length, while dense retrieval methods exhibit an inverse relationship. This finding suggests that there is potential for developing more sophisticated ensemble methodologies that could fully leverage the strengths of both approaches.

G Prompting Approach

In our experiments, we adopt the prompting approach of *Locate-and-Answer*, to facilitate active acquisition of information in the context when performing question answering. We first request LLMs to identify and locate the corresponding sentence that include the knowledge for the question from the provided document, and then reason to provide its answer. According to the experimental result in Table 3, this prompting approach can effectively improve the QA performance of LLMs compared to direct answering.

H Calibration-Induced Performance Discrepancies

In our study, retrieved documents serve as the unequivocal reference, mirroring practices in K-12 education. However, the evaluation results reveal counterintuitive calibration patterns in state-of-theart LLMs, particularly regarding their performance on *simple direct* questions versus more complex question types. For instance, stronger models such as GPT-4 and GPT-40 may underperform weaker ones on *simple direct* questions, which primarily assess factual recall. This discrepancy suggests that, for *simple direct* questions, the calibration of LLMs may lead models to exhibit excessive confidence in their internal knowledge, thereby inhibiting their reliance on external documents as the ground truth.

In real-world situations, however, knowledge conflicts typically occur in less established factual domains. As a result, stronger LLMs are expected to demonstrate improved performance in such scenarios due to their ability to better navigate and resolve these conflicts.

Original Paragraph	The halophiles, which means "salt-loving", live in environments with high levels of salt. They have been identified in the Great Salt Lake in Utah and in the Dead Sea between Israel and Jordan, which have salt concentrations several times that of the oceans				
Updated Paragraph	 The halophiles, which means "pressure-loving", live in environments with high levels of pressure. They have been identified in the Mariana Trench in western Pacific Ocean, which have higher pressure than other environments 				
Question	In which type of environments do halophiles typically live?A. High Acidity Environments.B. High Salt Environments.C. High Pressure Environments.D. High Sugar Environments.				
Original Answer	B. High Salt Environments.				
Updated Answer	C. High Pressure Environments.				

Table 5: An example of hypothetical knowledge update for a question in Biology. The modifications of factual knowledge and contextual information in the paragraph are highlighted in *red*.

Document Statistics		Total				
	Chem.	Bio.	Phys.	Geo.	Hist.	
Average Document Length Num of Documents	315.5 291	588.1 671	647.3 166	429.8 471	277.1 606	437.3 2205

Table 6: Average length (in word counts) and quantities of documents in different subjects.

# Ouestions (avg. length)	Question Types								
··· C (B)	Simple Direct	Multi-hop Direct	Multi-hop Distant	Multi-hop Implicit	Distant Implicit				
Chemistry	73 (12.4)	88 (17.8)	78 (17.9)	40 (17.4)	68 (21.1)	347 (17.3)			
Biology	148 (11.1)	170 (16.7)	248 (16.7)	94 (16.7)	213 (20.7)	873 (16.7)			
Physics	46 (12.0)	49 (16.7)	49 (15.6)	41 (16.9)	39 (19.9)	224 (16.1)			
Geography	141 (12.1)	147 (16.7)	162 (16.8)	96 (17.7)	144 (21.3)	690 (16.9)			
History	202 (12.8)	201 (16.1)	171 (17.0)	146 (18.1)	151 (21.6)	871 (16.8)			
Total	610 (12.1)	655 (16.7)	708 (16.9)	417 (17.5)	615 (21.0)	3005 (16.8)			

Table 7: Average question length (in word counts) and quantities of questions in different subjects and types.

Retrieval Methods	Subjects				Question Types					Average	
	chem.	bio.	phys.	geo.	hist.	Sim. Dir.	Mul. Dir.	Mul. Dis.	Mul. Imp.	Dis. Imp.	0
BM25	87.90	84.65	74.55	78.26	81.52	79.34	85.04	86.30	75.78	84.23	82.73
Mistral-embed Ada-002	84.73 84.44	77.09 79.50	83.48 82.14	73.48 72.75	80.94 81.29	82.13 82.30	79.08 80.76	79.38 80.08	80.81 80.10	72.85 73.01	78.74 79.23

Table 8: Retrieval performance (in recall@1 %) of retrieval approaches in different subjects and question types.

Large Language Models			Subjects			Question Types					Average
	chem.	bio.	phys.	geo.	hist.	Sim. Dir.	Mul. Dir.	Mul. Dis.	Mul. Imp.	Dis. Imp.	
Llama3-8b GPT-3.5-turbo GPT-4-turbo	77.81 78.93 85.30	75.60 80.53 89.23	75.45 82.59 86.16	80.14 83.04 90.00	78.65 82.43 90.13	90.33 92.62 95.74	88.85 90.53 96.18	89.69 91.24 96.19	63.55 71.22 81.06	49.43 57.72 71.71	77.77 81.73 88.99

Table 9: Performance (in accuracy %) of LLMs in question answering in different subjects and question types.

_

Model	Prompting Method					
	Direct Answer	Locate-and-Answer				
Gemini-1.5-flash GPT-3.5-turbo GPT-40	75.21 73.91 80.43	78.54 81.73 87.09				

Table 10: Performance (in accuracy %) of LLMs in question answering with different prompting methods.