Knowledge Representation Approaches for Educational Question Generation

Anonymous ACL submission

Abstract

Teachers are increasingly using prompted LLMs to generate exam questions, and students can use generated questions for self-assessment. When generating questions from a given educational text-rather than relying solely on the LLM's internal knowledge—handling long textual content, such as a textbook spanning hundreds of pages, presents a challenge. In this paper, we experiment with three knowledge representation approaches tailored for educational question generation using LLMs. As a novel contribution among these alternatives, 013 we adapt the atomic fact decomposition method from fact-checking research to the educational domain. We manually evaluate the generated questions based on various criteria. Our empir-017 ical results indicate that a list of atomic facts provides a better foundation for question generation than long plain text and that LLM-based 019 question generation from Knowledge Graph 021 triplets outperforms rule-based question generation from Knowledge Graphs.

1 Introduction

024

The automatic generation of educational questions (EQG) will play a key role in scaling education (Baidoo-Anu and Ansah, 2023). It can significantly decrease the workload of teachers along with enabling student self-assessment at scale in a personalized way. Large language models (LLMs) have demonstrated great performance in various tasks and applications. Teachers also employ LLMs to generate assessment questions for their exams by prompting general models (Walter, 2024).

In this paper, we focus on EQG given textual educational material (textbooks, lecture slides, transcripts, etc.). Questions are generated faithfully to the educational material and not from the latent knowledge of pre-trained general LLMs. This use case is important in rapidly evolving fields, as it can be applied even to web content like blog posts. Onboarding processes in organizations can also exploit automated assessments against their own document base. Lastly, a teacher might want to teach special topics that should be included in assessment questions. We note that this approach seems to be very similar to RAG, but in EQG there is no query, thus instead of the retriever step it requires special techniques to gather relevant contexts from the educational material for the generation step. 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Although LLMs have great potential in EQG, their use is not unproblematic, as they suffer from hallucinations and misinformation. Chen et al. (2024) and our preliminary experiments also report that LLMs can generate close to perfect simple questions. Simple questions are remember-type questions (according to the Bloom's Taxonomy of educational science (Anderson et al., 2001)) which can be generated from a short chunk of text as relevant content. Increasing the complexity level of questions requires the LLM to get to a deeper semantic understanding of the text; thus, it leads to hallucinations. Similarly EQG from long contents, for instance a textbook, also introduces more mistakes as distant relationships in text are more difficult to detect in LLMs. In this paper, we propose four methods for generating multiple-choice questions (MCQs) from multiple chapters of textbooks. These types of items are still remember-type questions, but the generation of good-quality questions and answer options requires comprehensive memorization of all of the educational material.

There are various opportunities to represent the knowledge of a larger educational material that provides the knowledge context for the generation step of the LLM. Knowledge Graphs (KGs) are structured representation formats that represent knowledge in entity-relation-entity triplets and can store unlimited amounts of information. The key practical disadvantage of KGs is that they are rigid, and building a KG from text is still not precise enough with moderate coverage (despite the vast amount

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

132

of research on KG extraction (Kertkeidkachorn and Ichise, 2017; Koncel-Kedziorski et al., 2019; Melnyk et al., 2022)). On the other end, we can consider the original long-form texts as a knowledge representation format, having perfect precision and coverage. The problem is that LLMs are not accurate enough to understand deeper and long-distance semantic/pragmatic relationships in long texts. We propose to use a solution that lies between KGs and long texts. Inspired by fact-checking research (Min et al., 2023; Chung et al., 2025), we extract atomic facts from the original text and use the set of atomic facts as the knowledge representation format for EQG. The definition of atomic fact is EQG-driven, it can be any simple and clear statement, i.e. whose knowledge can be assessed. Our assumption is that the set of atomic facts, which are simple and clear sentences, are considerably easier to understand by LLMs than long texts and provide a less rigid and better coverage representation than KGs could.

083

087

100

101

102

105

106

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

We introduce four EQG methods grounded on the three knowledge representation formats (two methods for the KGs). Each of the methods employs GPT-40 (Hurst et al., 2024) with few-shot prompting. We carried out comparative experiments on the methods and generated comparetype MCQs for sorting algorithms. The generated MCQs were carefully evaluated by human experts (teachers of algorithms subject) according to seven evaluation metrics. The evaluation metrics were designed for real-world usability, i.e. whether they are usable by a teacher or for student self-assessment. The metrics cover areas from factual correctness through clarity to creativity and engagement. The conclusion of these experiments is that there is no clear winner, but there is a considerable tradeoff between question correctness.

In summary, our contributions are threefold in this paper:

- We adapt the atomic fact text decomposition method to the EQG task.
- We propose and comperatively evaluate four methods for educational MCQ generation based on three different types of knowledge representation.
- Our human evaluation methodology for generated MCQs consists of various aspects that are important in their real-world usability, including factuality, clarity, creativity, and engagement metrics.

2 Related Work

Multiple-choice question generation Multiplechoice question (MCQ) generation has been extensively studied in educational technology.

Early research focused on template-based or rulebased techniques. Papasalouros et al. (2008) developed algorithms for MCQ generation based on domain ontologies, introducing eleven strategies. Of these eleven, Strategy 6 and Strategy 7 are very similar to our rule-based KG approach.

LLMs have been disrupting the EQG field in the last two years. For instance, Maity et al. (2024) introduced multi-stage prompting with GPT models, showing improved performance across multiple languages. Scaria et al. (2024) reported that LLMs are able to generate MCQs keeping the Bloom's Taxonomy levels, while Yao et al. (2024) and Wang et al. (2024) developed a self-refine framework for professional exam questions using iterative selfcritique.

Recent work demonstrates, that hybrid methods are still relevant. Kumar et al. (2024) combined ontology-based and machine learning techniques for MCQ generation. Their approach also specifically addresses generating questions for different cognitive levels.

Prior EQG work focuses on generating questions either from localized contexts or without any grounding in a zero-shot manner. Our work addresses the challenge of generating questions that require information scattered across multiple chapters or sections of educational material.

Atomic fact extraction Atomic fact extraction is a recently popularized technique mainly used in claim-verification. Min et al. (2023) argues that complex claims could contain both valid and invalid information, hence the need for decomposing these complex sentences into statements that each contain a single piece of information, that is unambiguous in their factuality.

Min et al. (2023) proposes a complete framework for evaluating LLM-generated text by breaking it down to atomic statements and estimating the ratio of such statements that were supported by a trusted source. Chung et al. (2025) presents a complex system, leveraging a similar approach in the medical domain to verify LLM-generated text based on the patient's medical history. Wanner et al. (2024) compares different methods for breaking down claims into atomic facts and shows that the decomposition method significantly affects

185

186

188

190

191

192

193

194

195

196

197

198

199

201

204

209

210

211

212

214

215

216

217

219

221

229

232

downstream results on factual precision measurements, such as FACTSCORE (Min et al., 2023).

Atomic fact extraction is useful in other domains besides claim verification. Chen et al. (2023) found that breaking down complex sentences into selfcontained factual statements, which they refer to as "propositions" significantly improves retrieval performance in a dense retrieval setting. They work with various open-domain QA datasets and English Wikipedia text, presenting FACTOIDWIKI, English Wikipedia broken down into various levels of granularity (passage, sentence, and proposition). Kamoi et al. (2023) uses GPT-3.5 to automatically decompose Wikipedia claims into subclaims, showing improved performance on entailment tasks across multiple datasets.

To the best of our knowledge, we are the first to employ atomic facts as educational knowledge representation. Our definition of atomic facts is slightly different from the claim-checking definition as we focus on small self-contained knowledge statements that are suitable for evaluation. The educational use case requires that the atomic facts be both verifiable and pedagogically relevant.

3 Question Generation Approaches

3.1 MCQ for Comparing Skills

We automatically generate MCQs, which is a common assignment format in education, as they can test various levels of complexity and cognition of students (Masters et al., 2001; Brady, 2005; Scaria et al., 2024). We address monodisciplinary EQG (Chen et al., 2024) and in this study, we experiment with comparison questions for sorting algorithms. To answer these questions, the test-taker has to remember the concepts in algorithm studies (e.g. time complexity), concepts that are specific to sorting algorithms (e.g. stable sorting), and has to remember the sorting algorithms' properties. According to Bloom's Taxonomy (Anderson et al., 2001), these questions belong to the remember level, as these MCQs can be answered with perfect remembering, while comprehensive understanding is not necessary.

On the other hand, remembering many pieces of information is necessary to answer the comparisontype MCQs, and this information is scattered in the educational material, requiring the test-taker to keep in mind dozens of pages. The key difference between this paper and related EQG work (Chen et al., 2024; Elkins et al., 2023) is that we generate questions from very long texts of educational material and a single MCQ can ask for facts mentioned in the texts far from each other. Precise and full coverage understanding of long texts is still a challenge for LLMs. We describe three knowledge representation forms for EQG, designed to overcome this issue.

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

277

278

279

280

3.2 Knowledge Representation Alternatives

We explored three approaches to represent the educational knowledge. A comparative example of the alternatives is presented in Figure 1.

Plain Text As LLMs have the ability to make sense of plain textual information, the first natural choice was to represent the knowledge as text. Our first approach utilized text descriptions, sourced from Wikipedia entries. In educational deployment, this could be extended to textbooks, lecture notes, and other trusted materials. This format maintains natural language context but requires LLMs to identify and extract relevant information.

Knowledge Graph KGs offer a structured representation through entities and relationships. The structure of KGs makes programmatic manipulation possible while meaningful edge types and concept nodes could make them interpretable for humans. A large number of KGs are available, and constructing new ones for specialized domains is also manageable.

It is important to note that KGs and knowledge spaces (Doignon and Falmagne, 1985, 2015), which are often used in the education literature, are separate approaches to organizing educational content. While knowledge spaces provide manageable chunks for direct instruction and assessment and focus on prerequisites, KGs are a general representation of data, making use of concepts and the relationships between them.

Factual Statements This approach transforms natural text into atomic fact statements, where each statement captures a single, self-contained piece of information.

Each factual statement must satisfy three criteria:

- 1. Non-triviality: The statement should convey meaningful domain knowledge
- 2. Self-containment: The statement should be comprehensible without additional context
- 3. Verifiability: The statement should be clearly true or false within the domain



Figure 1: Overview of four EQG methods over three knowledge representation forms.

This representation combines the accessibility of natural language with some of the structural benefits of KGs, as statements often follow implicit subject-predicate-object patterns.

3.3 MCQ Generation Methods

284

286

287

290

291

299

303

304

Building upon the previously described knowledge representations, we developed four distinct methods for generating multiple-choice questions, see Figure 1 for an overview. Each method uses LLMs and we engineered the EQG prompts to be as similar to each other as possible.

Plain Text-Based Method A straightforward approach would be to give the entire educational material to the LLM to generate questions. However, due to the limited context window of LLMs, this approach quickly becomes infeasible with large volumes of content. Splitting the data is not an option either, as it would prevent generating questions that consider the entire document. To address these challenges, we partitioned the source material into manageable segments (e.g. textbook chapters) and first applied a summarization process to each segment individually. This process removes less important details while still allowing the LLM to grasp the big picture. The resulting summaries were then used as input for the few-shot question generation prompt.

308Rule-based EQG from Knowledge GraphsTo309exploit the structured information inherent in310KGs, we created a rule-based question generation311methodology that utilizes predefined templates cor-312responding to the graph's nodes and edges. Two313primary templates were constructed. In the first

template, the inquiry focuses on identifying which 314 algorithm exhibits a specified property. Let knowl-315 edge graph $G = \{(h, r, t) | h, t \in V, r \in R\}$ is a 316 set of triplets, where V is the set nodes and R is 317 the set of relation types. Let $p \in V$ a selected 318 property and $r \in R$ a selected relation type, and 319 we have four sorting algorithms $s_a \in S$ the answer 320 and $s_{1-3} \in S$ distractors where $S \subset V$. 321

Which algorithm has $\{r\}$ of $\{p\}$?	322
- s_a , where $(s_a, r, p) \in G$	323
- s_1 , where $(s_1, r, p) \notin G$	324
- s_2 , where $(s_2, r, p) \notin G$	325
- s_3 , where $(s_3, r, p) \notin G$	326

327

328

330

336

337

338

339

340

341

343

344

In the second template, the focus is on a specific algorithm $s \in S$ by inquiring about one of its properties. In this case, the question is formulated as:

What is the $\{r\}$ of $\{s\}$?	331
- p_a , where $(s, r, p_a) \in G$	332
- p_1 , where $(s, r, p_1) \notin G$	333
- p_2 , where $(s, r, p_2) \notin G$	334
- p_3 , where $(s, r, p_3) \notin G$	335

Where $p_a \in V$ the answer and the $p_{1-3} \in V$ distractors.

Because the generated questions initially relied exclusively on the formal structure of the KG, this approach led to many unnatural questions for example, *"Which algorithm has instance of of adaptive sort?"*. We prompt an LLM to rephrase these questions to make them more understandable for students. 345Knowledge Graph + LLMIn this method, in-346stead of relying on predefined templates, we use the347KG – or a selected subset thereof – as the context of348an LLM. The KG is represented in a list of textual349(h, r, s) triplets. The LLM few-shot prompted for350MCQ generation based on this context can infer351the relationships necessary for coherent question352generation.

Atomic Facts The last method employs atomic facts as context. These factual statements have a structure that closely resembles how knowledge is organized in a KG, while it is more flexible as there is no ontology schema. Atomic facts are more concise than presenting the same information as lengthy paragraphs of plain text. The list of atomic facts provides the context for the LLM employing a few-shot prompt very similarly to the previous method.

4 Experimental Results

4.1 Datasets

354

357

364

372

374

379

Our goal is to generate MCQs that require having information from different parts of educational materials, for instance, connecting concepts from different chapters of a textbook. Let S = Merge*sort*, *Selection sort*, *Heapsort*, *Timsort*, *Quicksort*, *Insertion sort*, *Bubble sort* represent selected sorting algorithms.

The set S of sorting algorithms represents a case of this general problem, where each algorithm can be considered a separate chapter within a textbook on computer algorithms. What makes sorting algorithms a good use case is that they share the same properties (time and space complexity, stability, etc.) while being different in important aspects. However, we believe that our approach can be generalized to most educational content. For each representation approach, we constructed specialized datasets:

Plain Text Corpus We extracted and processed
the English Wikipedia articles corresponding to
each algorithm in S. Original long textual content
was summarized by prompting an LLM, in order to
fit the most amount of text possible into the context
window of the LLM while maintaining essential
algorithmic concepts. The resulting corpus had a
total of 2,060 words.

Knowledge Graph Dataset We used Wikidata¹ as our source of KG, as it is naturally aligned with Wikipedia, thus trying to align our structured and unstructured knowledge representations. Due to the amount of information represented in the text, as opposed to any KG, this is not fully achievable. We extracted a subgraph from Wikidata centered on the algorithms in *S*. The graph was pruned to remove technical edges (e.g., entity IDs) while preserving meaningful relationships. The resulting graph contained 99 triplets.

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

Factual Statement Dataset We prompt an LLM to transform the plain text corpus into atomic factual statements. Each paragraph was sequentially processed by the model, which had access to the requirements and factual statement criteria. The final dataset comprised of 1,471 sentences. We note that the extraction of atomic facts is not a summarization or text simplification process. Instead, it is aimed at decomposing rather complex sentences into a pedagogically relevant set of simple, self-contained statements.

Topic-Oriented Filtering Topic-oriented filtering is an opportunity for the educator to customize the underlying dataset to contain only knowledge that the educator wants to assess. Filtering is the most convenient with KGs through removing triplets with certain relation types and atomic fact sentences can be simply deleted. On the other hand, modifying the source plain text to keep only pedagogically relevant content requires a lot of effort and uncertainty.

To align with educational objectives, we applied additional filtering across all representations to emphasize the technological aspects of sorting algorithms. This included removing historical development information, and focusing instead on operational characteristics and complexity analysis. To achieve this, in the case of Wikipedia, a summarization LLM was prompted to remove undesired content, reducing the corpus from 23,643 words to 2,060 words. The factual statements were reduced from 1,471 to 1,435 sentences using heuristical techniques to locate and remove this information. In the case of the KG, we collected certain relation types that indicate the presence of this type of information, such as WikidataProperty : P61("discoverer or inventor") and deleted all triples from the KG with these relation types, reducing the

¹https://www.wikidata.org/wiki/Wikidata

533

534

535

488

489

490

number of triplets from 99 to 85

440

441

442

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

482

483

485

486

487

For the sake of reproducibility, the cleaned plain texts, the KG triplets, and the list of atomic facts are available at github.com/anonym.

4.2 Experimental Details

We used GPT-40 (version 2024-08-06) with a temperature setting of 0.4 for all experiments, including summarization, atomic fact decomposition, and rephrasing tasks, accessing the model through OpenAI's API² from a local machine (a total API cost of \$10 USD). In preliminary experiments, we found that a lower temperature option (0.4) produced higher quality questions compared to a higher setting (0.7). No additional computational resources were required for the experiments.

The whole knowledge base fits into the context window of the prompts in all of the EQG prompts. The summarized Wikipedia articles consist of 2,060 words, while the atomic facts 18,021 words. In the Knowledge Graph + LLM method, the triplets were formatted as "*subject – predicate* $\rightarrow object$ ", allowing the model to access all relevant relationships during question generation.

At the KG rules method, when generating questions using the first template type, if a randomly selected property was not shared by four different algorithms (which was needed for one correct answer and three distractors), we randomly sampled algorithms as distractors, under the assumption that the absence of a property in the KG indicates the algorithm lacks that characteristic. This enabled us to generate questions about properties not shared by at least four algorithms in the dataset, increasing the pool of possible questions despite the number of triplets available.

All approaches used prompts that enforced: (1) inference across multiple facts for comparisonbased questions, (2) verification of answer choices by prompting to include evidence for both correct answers and distractors, and (3) grounding in explicitly stated information only. All prompts are provided in Appendix A.

4.3 Evaluation methodology

We applied seven evaluation metrics to compare the questions generated by different methods. These metrics were inspired by the evaluation methodologies of Chen et al. (2024); Elkins et al. (2023); Luo et al. (2024); Scaria et al. (2024) and were further refined based on insights gained during the annotation process.

Answer Correctness (Yes/No) Yes, if the provided answer is correct for the given question.

Distractor Incorrectness (Yes/No) Evaluates whether the distractors are valid, ensuring that none of the distractors constitute correct answers.

Option Quality (Good/Bad) Option quality is good if the correct answer and the distractors are non-trivial to distinguish and the MCQ is sufficiently challenging. Distractors should also avoid redundancy among them and trivial elimination. For example, the question should not implicitly reveal the answer, as in the case of "Which sorting algorithm uses a heap?"

Text Quality (1-5) Measures the linguistic quality of the question, assessing grammatical correctness, clarity, and readability. A score of 1-2 corresponds to grammatically incorrect sentences, a score of 3 indicates grammatically correct but not entirely clear questions, while a score of 4 represents understandable phrasing that could still be improved.

Compactness (Yes/No) Early experiments revealed generated MCQs with unnecessary details, such as asking about multiple concepts while the answer options can be chosen by knowing only one concept. This metric is yes if the question contains only essential information.

Template-Based vs. Creative (1-5) A very subjective metric, to assess the extent to which the question follows a template-based structure, i.e. boring (1) or diverse and exhibits creativity (5).

Would You Use It? (1-5) Evaluates whether a teacher would consider using the question in an exam. This is the most subjective metric, as it depends entirely on the preferences of the educator.

Prior studies have incorporated relevance-based metrics (Chen et al., 2024; Elkins et al., 2023; Luo et al., 2024; Scaria et al., 2024); however, we leave out this category because our preliminary findings indicated that all generated MCQs consistently met relevance criteria. Similarly, a frequently used metric evaluated whether a question aligned with a predefined category of the Bloom's Taxonomy (Bloom et al., 1956). Since our initial annotation process confirmed that all questions adhered to their respective categories, we excluded this metric as well.

²https://platform.openai.com/

	Answer Corr.	Distractor Incorr.	Compact.	Option Quality	Text quality	Template v. Creative	Would You Use It?
Plain Text	0.85	0.80	0.55	1.00	4.55	4.04	3.50
Atomic Facts	0.95	0.80	0.95	0.98	4.73	3.74	3.90
KG + LLM	0.90	0.95	1.00	0.95	4.71	2.65	3.78
KG Rules	0.90	0.75	1.00	0.79	4.70	2.51	3.36

Table 1: The four annotators' average scores on three factual correctness and four engagement metrics.

The evaluation of the generated MCQs was conducted by four annotators, all of whom are either current or former teachers of the Algorithm and Data Structures undergraduate course. Each annotator reviewed and annotated every example in the dataset.

4.4 Results

536

537

541

542

543

544

545

546

550

552

554

559

562

564

565

566

567

571

573

574

577

The performance of the methods was compared over an MCQ bank containing 20 generated questions per method (listed in Appendix B). Every MCQ was independently evaluated by all four annotators. To ensure unbiased evaluation, the MCQs were presented in a randomized order, and the annotators were unaware of their origin.

Answer Correctness, Distractor Incorrectness, and Compactness metrics are objective ones. The cases of disagreement among annotators were further examined by a designated annotator, who made the final decision. Table 1 presents the average scores of the annotators.

For the metrics where disagreement was not resolved, we conducted an analysis of inter-annotator agreement. In the case of the categorical *Option Quality* metric, we computed Fleiss' kappa (κ) across the four annotators, where the $\kappa = 0.89$ indicates a high level of agreement.

For the remaining nominal metrics, Spearman's rank correlation coefficient was employed to compare the decisions of the annotators. We report here the average of the pairwise correlations. Among the three nominal categories, the Template-Based vs. Creative metric achieved the highest correlation, with a coefficient of 0.61. In contrast, the Text Quality and Would You Use It? metrics exhibited lower agreement levels, with correlation coefficients of 0.31 and 0.22, respectively. The low correlation for Would You Use It? is understandable because this subjective measure is highly dependent on personal preferences. Text quality ratings received mostly 4 and 5 ratings. Annotators judged only nuances, thus the low agreement here can be attributed here to personal preferences also.

Our results reveal significant differences in MCQ factual correctness across the four approaches, with structured and semi-structured knowledge representations showing more favorable results. While the Plain Text approach achieved respectable scores, the aggregate performance across all factual metrics (*Answer Correctness, Distractor incorrectness* and *Option Quality*) shows that both Atomic Facts and KG + LLM achieved higher overall factual accuracy. This suggests that structured and semi-structured representations provide a more reliable grounding for EQG. On the other hand, the more freedom the method gets, the more creative MCQs are generated.

578

579

580

581

582

583

584

585

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

5 Discussion

5.1 Error analysis

The Would You Use It? metric indicates a preference for MCQs generated using either atomic facts or KG triplets as LLM context. This can be partly due to the conciseness of the MCQs generated by these approaches, illustrated by the *Compactness* score. The plain text-based approach produced many questions with redundant content, this being the reason for its low score on the *Compactness* metric. For instance, the question "Which sorting algorithm is described as a stable, hybrid sorting algorithm that combines merge sort and insertion sort?" provides multiple attributes when a smaller subset of information would uniquely identify Timsort as the correct answer for the question.

The rule-based KG approach, despite having the least amount of LLM influence, did not achieve the high factual accuracy we expected. KG incompleteness is one source of these errors, but rephrasing can also introduce mistakes due to ambiguities in edge labels. For instance, the ambiguity of the triplet "*insertion sort – derivative work* \rightarrow *Timsort*" led to the following rephrased question: "Which algorithm is a derivative of Timsort?" where the direction of the "*derivative*" is changed.

Another common issue was the rigidity of the rule-based approach, which struggled with incon-

sistencies in the KG. For instance, the interchangeable use of "subclass of" and "instance of" relationships in Wikidata ("bubble sort – instance of" \rightarrow sorting algorithm" vs. "merge sort – subclass of" \rightarrow sorting algorithm") created problems for rule-based generation. The KG + LLM approach overcame these issues by leveraging the LLM's ability to understand semantic similarities between edge types and node names, leading to better performance in distractor validity.

5.2 Hallucination

620

625

626

633

635

637

641

642

649

651

652

654

657

670

We address EQG's faithfulness to the given educational material. As the topic of computer algorithms is well known, the basic concepts are probably learnt by the LLMs. To measure the methods' hallucination, we designed a standalone experiment in which deliberately false information was injected into the data sources. The objective was to determine whether the generated MCQs would rely on the falsified data or the pretrained knowledge of the LLM.

Initially, we permuted the names of the sorting algorithms present in the preprocessed texts (summarized texts and atomic facts) and in the KG. Specifically, we applied a permutation where every occurrence of the algorithms' names was replaced, carefully addressing variations such as "*mergesort*" versus "*Merge sort*". The below mapping is applied to the original algorithm identifiers:

selection sort \rightarrow heapsort \rightarrow insertion sort \rightarrow \rightarrow Timsort \rightarrow bubble sort \rightarrow \rightarrow quicksort \rightarrow merge sort \rightarrow selection sort

After the replacement, the four EQG methods were executed on the permuted dataset, yielding new MCQs derived from the intentionally corrupted data. To ensure that the generated MCQs remained comparable, we subsequently revered the replacement over the new MCQs and options to restore the original algorithm names. The resulting MCQ bank was then subjected to an annotation process using only the *Answer Correctness* and *Distractor Incorrectness* metrics. For each method, 20 MCQs were evaluated and the outcomes for the corrupted MCQs are presented in Table 2.

The results show that the KG Rules method, which depends least on the capabilities of the LLM, exhibited the lowest degree of hallucination. In this method, the LLM's primary task was to rephrase the MCQ, ensuring that the answers remained consistently aligned with the KG. The mistakes of KG

	Answe	er Corr.	Distractor Incorr.		
	Original	Corrupted	Original	Corrupted	
Plain Text	0.85	0.60	0.80	0.50	
A. Facts	0.95	0.60	0.80	0.20	
KG + LLM	0.90	0.65	0.95	0.20	
KG Rules	0.90	0.95	0.75	0.60	

Table 2: Hallucination experiments: Scores on original and corrupted datasets.

Rules in the corrupted data has the same source as the original data, i.e. relation direction ambiguity. In contrast, for the other three methods, wherein the LLM was responsible for constructing the MCQ, the incidence of errors was much higher and the three methods suffered from halucination at the same level. 671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

The Plain Text method achieved considerably better scores on *Distractor Incorrectness* compared to other methods. Analysing the MCQs, this can mainly attributed to behaviour of the method generating distractors that were specific to algorithmic contexts, such as "Using a heap data structure."

6 Conclusions

In this paper, we analysed the knowledge representation formats of long educational textual materials to automatically generate comparison-type MCQs. Our experiments show that generating EQG directly from long, plain texts gives significantly more factually incorrect answers than other representations, while rule-based MCQ generation from KG yields boring MCQs and question quality highly depends on the coverage and completeness of the KG. Instead of these methods, we recommend to list facts as the context of an EQG prompt and exploit LLMs to generate questions. The list of facts can be triplets from a KG or simple statement sentences (atomic facts) decomposed from the original long, plain text. The choice between these methods should consider two factors: (1) the availability of a high-quality KG for the domain, and (2) the desired trade-off between factual accuracy and engagement. While both approaches demonstrate overall good factual correctness, KG + LLM scoring slightly higher, the atomic facts method produces more creative and engaging questions, as indicated by its higher "Template-Based vs. Creative" metric compared to KG + LLM.

709 Limitations

711

712

713

714

715

716

717

718

719

720

721

724

725

731

733

736

737

740

741

742

743

745

746

747

In this study, we experimented exclusively with a narrow domain of sorting algorithms. A crucial limitation of our findings is that we do not know whether they hold on characteristically different other domains, like arts.

Similarly, we only generate comparison-type MCQs. We can only guess that other types of questions requiring remembering distant facts would behave in the same way.

We employed hand-crafted and manually cleaned KG as a knowledge representation option. In the real world, KG has to be constructed from text and the text-to-KG process is far from perfect (see the review of Pan et al. (2024)). The errors introduced by the text-to-KG process might decrease the efficiency of the KG+LLM method.

In this study, we assumed that the subset of the educational material – either in the format of plain text or a list of atomic facts or KG triplets – from which comparison-type MCQs can be generated is given. This is not available in real-world situations where you can have hundreds of textbook pages. As a future work, we are proposing methods to identify a reduced number of subsets of distant text chapters or subsets of atomic facts which gives an appropriate grounding for EQG.

Our evaluation setup is human labor intensive. As a future work, it would be interesting to investigate whether some of our metrics can be replaced by automatic evaluation in an LLM as a judge approach. The annotations of the four domain experts for the 160 generated MCQs provide a useful benchmark for the comparative evaluation of human and LLM as a judge evaluations.

As Section 5.2 shows, hallucination is a serious issue with the proposed methods. Improving faithfulness is one of our most important research challenges.

48 Acknowledgments

anonym

749

References

LW Anderson, DR Krathwohl, PW Airasian, PA Cruikshank, Richard Mayer, RJ Pintrich, J. Raths, and MC Wittrock. 2001. A Taxonomy for Learning, Teaching, and Assessing: A Revision of Bloom's Taxonomy of Educational Objectives. 750

751

752

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

- David Baidoo-Anu and Leticia Owusu Ansah. 2023. Education in the era of generative artificial intelligence (ai): Understanding the potential benefits of chatgpt in promoting teaching and learning. *Journal of AI*, 7(1):52–62.
- B. S. Bloom, M. B. Engelhart, E. J. Furst, W. H. Hill, and D. R. Krathwohl. 1956. *Taxonomy of educational objectives. The classification of educational goals. Handbook 1: Cognitive domain.* Longmans Green, New York.
- Anne-Marie Brady. 2005. Assessment of learning with multiple-choice questions. *Nurse Education in Prac-tice*, 5(4):238–242.
- Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang, and Dong Yu. 2023. Dense x retrieval: What retrieval granularity should we use? *arXiv preprint arXiv:2312.06648*.
- Yuyan Chen, Chenwei Wu, Songzhou Yan, Panjun Liu, and Yanghua Xiao. 2024. Dr. academy: A benchmark for evaluating questioning capability in education for large language models. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3138–3167.
- Philip Chung, Akshay Swaminathan, Alex J Goodell, Yeasul Kim, S Momsen Reincke, Lichy Han, Ben Deverett, Mohammad Amin Sadeghi, Abdel-Badih Ariss, Marc Ghanem, et al. 2025. Verifact: Verifying facts in Ilm-generated clinical text with electronic health records. arXiv preprint arXiv:2501.16672.
- Jean-Paul Doignon and Jean-Claude Falmagne. 1985. Spaces for the assessment of knowledge. *International journal of man-machine studies*, 23(2):175– 196.
- Jean-Paul Doignon and Jean-Claude Falmagne. 2015. Knowledge spaces and learning spaces. *arXiv* preprint arXiv:1511.06757.
- Sabina Elkins, Ekaterina Kochmar, Iulian Serban, and Jackie CK Cheung. 2023. How useful are educational questions generated by large language models? In *International Conference on Artificial Intelligence in Education*, pages 536–542. Springer.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and

Greg Durrett. 2023. Wice: Real-world entailment for

claims in wikipedia. In Proceedings of the 2023 Con-

ference on Empirical Methods in Natural Language

Natthawut Kertkeidkachorn and Ryutaro Ichise. 2017.

T2kg: An end-to-end system for creating knowledge graph from unstructured text. In Workshops at the

Thirty-First AAAI Conference on Artificial Intelli-

Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan,

Archana Praveen Kumar, Ashalatha Nayak, Man-

jula Shenoy K, Chaitanya, and Kaustav Ghosh. 2024.

A novel framework for the generation of multiple

choice question stems using semantic and machinelearning techniques. International Journal of Artifi-

cial Intelligence in Education, 34(2):332–375.

Haohao Luo, Yang Deng, Ying Shen, See Kiong Ng, and

Tat-Seng Chua. 2024. Chain-of-exemplar: Enhanc-

ing distractor generation for multimodal educational

question generation. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7978-

Subhankar Maity, Aniket Deroy, and Sudeshna Sarkar.

Joan C Masters, Barbara S Hulsmeyer, Mary E Pike, Kathy Leichty, Margaret T Miller, and Amy L Verst.

2001. Assessment of multiple-choice questions in

selected test banks accompanying text books used in

Igor Melnyk, Pierre Dognin, and Payel Das. 2022.

Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision

in long form text generation. pages 12076-12100.

Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap.

IEEE Transactions on Knowledge and Data Engi-

Andreas Papasalouros, Konstantinos Kanaris, and Konstantinos Kotis. 2008. Automatic generation of mul-

Knowledge graph generation from text. In Find-

ings of the Association for Computational Linguistics:

2024. A novel multi-stage prompting approach for

language agnostic mcq generation using gpt. In European Conference on Information Retrieval, pages

formers. arXiv preprint arXiv:1904.02342.

Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text

generation from knowledge graphs with graph trans-

Processing, pages 7561–7583.

gence.

7993.

268–277. Springer.

nursing education.

EMNLP 2022, pages 1610-1622.

- 808
- 810 811
- 812
- 814
- 815 816
- 817 818
- 819 820
- 821 822

- 825 826
- 830
- 831
- 832 833
- 834 837
- 838
- 841 842

844

- 849

852

853 854

tiple choice questions from domain ontologies. e-Learning, 1:427-434.

neering.

Nicy Scaria, Suma Dharani Chenna, and Deepak Subramani. 2024. Automated educational question generation at different bloom's skill levels using large language models: Strategies and evaluation. In International Conference on Artificial Intelligence in Education, pages 165–179. Springer.

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

- Yoshija Walter. 2024. Embracing the future of artificial intelligence in the classroom: the relevance of ai literacy, prompt engineering, and critical thinking in modern education. International Journal of Educational Technology in Higher Education, 21(1):15.
- Jiayi Wang, Ruiwei Xiao, and Ying-Jui Tseng. 2024. Generating ai literacy mcqs: A multi-agent llm approach. arXiv e-prints, pages arXiv-2412.
- Miriam Wanner, Seth Ebner, Zhengping Jiang, Mark Dredze, and Benjamin Van Durme. 2024. Α closer look at claim decomposition. arXiv preprint arXiv:2403.11903.
- Zonghai Yao, Aditya Parashar, Huixue Zhou, Won Seok Jang, Feiyun Ouyang, Zhichao Yang, and Hong Yu. 2024. Mcqg-srefine: Multiple choice question generation and evaluation with iterative self-critique, correction, and comparison feedback. arXiv preprint arXiv:2410.13191.

A	Prompts	883
This	s section shows the prompts that were applied to the results of the paper.	884
A.1	Plain Text-Based Method	885
The	following prompt was used to summarize Wikipedia articles.	886
1	Remove any unnecessary text from the given passage, retaining only the essential information needed to understand the main \rightarrow topic, while ensuring individuals familiar with the topic gain the same insights as before.	

3	# Steps
4	
5	1. **Read the Text** : Carefully review the entire passage to comprehend its main topic and context.
6	2. **Identify Key Points** : Determine the primary ideas, statements, or data that are crucial for understanding the topic,
	\hookrightarrow ensuring all necessary insights remain for those familiar with the subject.
7	3. **Filter Out Extraneous Information** : Remove any superfluous text such as filler words, repetitive phrases, tangential
	\hookrightarrow details, and unrelated anecdotes that do not contribute to the core understanding of the topic.
8	4. **Maintain Coherence**: Ensure the remaining text is clear and coherent without compromising the primary message or
	\hookrightarrow understanding of the topic, retaining essential insights for knowledgeable readers.
9	
10	# Output Format
11	
12	A concise and focused text that contains only the necessary information needed to understand the main topic, ensuring
	\hookrightarrow comprehensive understanding by all audiences.
13	
14	# Notes
15	

- 16 Consider the audience for which the text is being tailored to ensure that all retained information is relevant and maintains \hookrightarrow the knowledge depth suitable for those familiar with the topic.
- $^{-}$ Be mindful of retaining context and logical flow even as you remove extraneous text, ensuring experts continue to learn at \hookrightarrow the same level of depth.

The following prompt was applied to generate MCQs from the summarized texts.

2

```
1
     Create multiple choice questions from the provided documents frathat assess Understanding" according to Anderson's taxonomy.
2
3
     Focus on generating questions that require learners to infer, and compare content from the documents.
4
     # Steps
5
6
7
     1. **Analyze the Documents:** Thoroughly read the attached documents to grasp key concepts and details necessary for creating
         insightful questions.
 8
     2. **Determine Key Concepts:** Identify the main ideas and important supporting details.
 9
     3. **Generate Questions:** Formulate clear and easy to understand questions that require the application of understanding
     \hookrightarrow skills, such as interpretation or summarization.
     4. **Create Options:** Develop plausible distractors for each question along with the correct answer.
10
     5. **Provide Context:** Include a context from the sources for each question to guide learners, based on the attached
11
     \hookrightarrow documents.
12
13
     # Output Format
14
     The output should be formatted as a JSON array of objects. Each object should represent a question and include:
15
     - **question**: The text of the question.
16
17
     - **options**: An array of answer choices.
18
     - **correct_answer**: The position or index of the correct answer in the options array. Starting from 0.
     - **grounding**: A context from the sources for the question, based on the attached documents.
19
20
     ···json
21
22
     Ε
23
       {
24
          "question": "What is the main concept addressed in [Section/Paragraph]?",
          "options": ["Option A", "Option B", "Option C", "Option D"],
"correct_answer": 1,
25
26
          "grounding": "Provide a context from the sources for the question."
27
28
       },
    ]
29
30
31
32
     # Examples
33
     **Example 1:**
34
35
36
     _Input:
37
     An article about an algorithm.
38
     An article about an another algorithm.
39
40
     _Output:_
41
        `json
     Ε
42
43
       {
         "question": "Which algorithm has instance of of divide-and-conquer algorithm?",
44
```

```
"options": ["heapsort", "quicksort", "Timsort", "selection sort"],
45
           "correct answer": 1.
46
           "grounding": "The divide-and-conquer algorithm is a key concept in quicksort."
47
48
       }
49
     ]
50
51
     (Real examples should be longer and customized to the attached documents, focusing on understanding concepts presented.)
52
53
54
     # Notes
55
56
     – Ensure each question is designed to test understanding within Anderson's taxonomy.
     Verify that distractors are plausible to encourage critical thinking.Tailor questions to reflect the specific content and style of the attached documents.
57
58
59
     - Provide context from the sources to guide learners in answering the questions.
```

Given a knowledge base of factual statements about a topic, generate 10 multiple choice questions that test understanding of

A.2 Atomic Facts

1

To generate questions from the atomic facts we used the following prompt:

the concepts and relationships described in these facts.

2 Focus on generating questions that require learners to infer, and compare content from the triplets. 3 4 # Steps 5 6 7 1. First, ****identify clusters of related facts**** by either: a) Finding predicates that connect multiple subjects to the same type of object (e.g., all algorithms with their space complexity) 8 9 10 OR b) Finding predicates that connect one subject to multiple objects 11 (e.g., all properties of one algorithm) 12 13 14 2. ****Before formulating a question, verify****: - You have enough distinct options for 4 meaningful choices 15 The relationships are unambiguous (one clear correct answer) 16 - Supporting evidence exists for both the correct answer and distractors 17 18 - For property questions: at least 4 different algorithms have this property - For value questions: at least 4 different possible values exist across algorithms 19 20 21 # Guidelines for questions 22 23 1. Answers must be definitively provable from the given facts 24 2. All 4 options must be plausible and related to the topic 25 3. Incorrect options (distractors) should be based on facts about other related concepts 26 4. The distractor answers should not be too trivial 27 5. Questions should test relationships between concepts or comparative properties 28 29 # Output Format 30 The output should be formatted as a JSON array of objects. Return nothing but the JSON response, pay attention to the format so 31 \hookrightarrow it can be loaded by Python's <code>`json.loads`</code> without any modifications. Each object should represent a question and include: 32 - **question**: The text of the question. - **choices**: An array of answer choices. 33 34 - **correct**: The position or index of the correct answer in the options array. Starting from 0. 35 - **grounding**: A list of factual statements from the sources for the question, based on the attached documents. 36 ···json 37 38 Ε 39 { "question": "...", 40 41 "choices": [$\frac{n}{n}$ \cdots $\frac{n}{n}$, 42 43 44 45 46 ٦. "correct": 0-3, // index of the correct answer "grounding": ["...", "...", "..."] // list the relevant facts supporting the correctness of the correct answer option 47 48 \hookrightarrow and any relevant information to the distractors 49 } 50] 51 52 53 # Examples 54 55 **Example 1:** 56 57 Input: List of factual statements about algorithms and their properties. 58 59 _Output:_ 60 61 ``json Ε 62

```
63
          {
               "question": "Which algorithm is not a stable sorting algorithm?",
64
               "choices": [
65
66
                   "Timsort"
67
                   "Heapsort",
                   "Insertion Sort",
68
                   "Bubble Sort"
69
              ],
"correct": 1,
70
71
72
               "grounding": [
73
74
75
                   "Timsort is a stable sorting algorithm",
                   "Heapsort is not stable",
"Insertion sort is a stable sorting algorithm",
76
                   "Bubble sort is a stable sorting algorithm"
77
              ]
78
          }
79
     ]
80
81
      **Example 2:**
82
83
84
      _Input:_
85
      List of factual statements about algorithms and their properties.
86
87
      _Output:_
88
        `json
      Ε
89
90
          {
91
               "question": "What is the worst-case space complexity of merge sort?",
               "choices": [
"0(n)",
"0(1)",
92
93
94
95
                   "O(n log n)",
96
                   "0(n^2)"
              ],
"correct": 0,
97
98
99
               "grounding": [
100
                   "Merge sort has a worst-case space complexity of O(n).",
                   "Merge sort has space complexity of O(n).",
101
                   "Quicksort has a worst-case space complexity of O(n).",
102
                   "Heapsort has a worst-case space complexity of O(1)."
103
104
              ]
105
          }
106
      ]
107
108
      # Rules
109
110
111
      1. Only use information explicitly stated in the fact statements
112
      2. Only reference concepts mentioned in the knowledge base
      3. Treat the facts as the sole source of truth
113
      4. If something isn't explicitly stated in a fact, don't assume it
114
115
      5. Consider different ways the same concept might be expressed
116
117
      # Notes
118
      - Verify that distractors are plausible to encourage critical thinking.
119
      - Provide triplets from the sources to guide learners in answering the questions.
120
121
122
      Return nothing but the json response, formatted to be loaded by Python's json.loads without any modifications.
```

The facts were extracted from the plain texts by the following prompt:

1	You are tasked with extracting fact-based statements from a text.
2	
3	Guidelines:
4	1. Break down complex sentences into simple, atomic facts. Each fact must clearly indicate its topic and be completely \hookrightarrow self-contained.
5	2. Only use information present in the source text
6	3. Only return pieces of information which are relevant for later testing knowledge on the topic in an educational setting (as \hookrightarrow the facts will be used for constructing test questions)
7	4. Focus on technical, definitional, and functional aspects that demonstrate understanding of the core concept
8	5. Each statement must clearly indicate what topic or concept it's describing, as if it could appear in a random set of facts
	\leftrightarrow about different topics
9	BAD: "Compares each element with the next one"
10	GOOD: "Bubble sort compares each element with the next element in the list"
11	
12	Input Format:
13	<input/>
14	Your text goes here
15	
16	
17	Output Format:
18	C
19	"Fact statement 1",

20	"Fact statement 2",
21	
22]
23	
24	Example:
25	
26	<input/>
27	The number ***#*** (/pai/; spelled out as "**pi**") is a mathematical constant, approximately equal to 3.14159, that is the
	\rightarrow ratio of a circle's circumference to its diameter. In addition to being irrational, $\star \pi \star$ is also a transcendental number,
	\rightarrow which means that it is not the solution of any non-constant polynomial equation with rational coefficients.
28	
29	
30	Output:
31	
32	" π is approximately equal to 3.14159",
33	" π equals the ratio of a circle's circumference to its diameter",
34	" π is an irrational number",
35	" π is a transcendental number".
36	"A transcendental number cannot be the solution of any non-constant polynomial equation with rational coefficients"
37	1
38	
39	Possible statements left out due to their lack of relevancy:
40	– "The number π is a mathematical constant"> as this is a very trivial statement

- " π is pronounced as 'pai'", --> this can't be effectively used when constructing a written test on π " π is spelled out as 'pi'", --> as this is a very trivial statement 41
- 42

A.3 Knowledge Graph + LLM

The below system prompt was applied to generate questions from the KG:

Given a knowledge base of facts in the form of triplets (subject--predicate->object), generate multiple choice questions that 1 \hookrightarrow test understanding of relationships and properties in the knowledge graph. 2

3 Focus on generating questions that require learners to infer, and compare content from the triplets.

4 5 # Steps

6 7

8

9

10

11

12 13 14

15

16

17 18 19

22

23

24

25

26 27

28 29

```
1. First, **identify clusters of related facts** by either:
```

- a) Finding predicates that connect multiple subjects to the same type of object
- (e.g., all algorithms with their space complexity)
- OR
 - b) Finding predicates that connect one subject to multiple objects
 - (e.g., all properties of one algorithm)

2. ****Before formulating a question, verify****:

- You have enough distinct options for 4 meaningful choices
- The relationships are unambiguous (one clear correct answer)

- Supporting evidence exists for both the correct answer and distractors
 For property questions: at least 4 different algorithms have this property
 For value questions: at least 4 different possible values exist across algorithms

20 21 3. **Document all relevant supporting triplets** that:

- Prove the correct answer
- Disprove incorrect options
- Use semantically similar predicates
- Validate the complete set of choices

Output Specifications

- 1. Generate 10 questions
- 30 2. Each question must have exactly 4 options
- 3. Only ONE option should be correct based on the knowledge base 31
- 32 4. All answers must be supported by explicit triplets
- 5. The distractor answers should not be too trivial 33 34

35 # Output Format 36

- The output should be formatted as a JSON array of objects. Return nothing but the JSON response, pay attention to the format so 37 ↔ it can be loaded by Python's json.loads without any modifications. Each object should represent a question and include:
- ****question****: The text of the question. 38 39 - **choices**: An array of answer choices
- **correct**: The position or index of the correct answer in the choices array. Starting from 0. 40
- ****supporting_triplets****: Supporting triplets from the input, keep the format. Do not modify the triplets. 41
- 42 ···json 43 44 Ε

```
45
         {
46
              "question": "...",
               "choices": [
47
                   ·...",
48
49
50
51
52
              "correct": 0-3, // index of the correct answer
53
```

```
"supporting triplets": [
54
55
                    // format: "subject--predicate->object, do not include any comments as it makes the result unparsable
               ]
56
57
          }
58
      ]
59
60
      # Examples
61
62
63
      **Example 1:**
64
65
       Input:
      Triplets from a knowledge graph, about algorithms and their properties.
66
67
68
      _Output:_
69
         `json
      Ε
70
71
           {
                "question": "Which algorithm is not a stable sorting algorithm?",
72
73
                "choices": [
74
                    "Timsort"
75
                    "Heapsort",
                    "Insertion Sort",
76
77
                     "Bubble Sort"
               ],
"correct": 1,
78
79
                "supporting_triplets": [
80
                     "Timsort--instance of->stable sorting algorithm",
81
82
                    "heapsort--has property->not stable",
                    "insertion sort--instance of->stable sorting algorithm",
"insertion sort--is a type of->stable sorting algorithm",
"bubble sort--instance of->stable sorting algorithm",
83
84
85
                    "bubble sort--is a type of->stable sorting algorithm",
86
                    "bubble sort--is->stable sorting algorithm'
87
               ]
88
89
           }
90
      ]
91
92
      Note how the same concept "being stable" is expressed through different predicates:
93
       - "instance of->stable sorting algorithm"
94
      - "is a type of->stable sorting algorithm"
95
      - "is->stable sorting algorithm"
96
      - "has property->not stable"
97
98
99
      **Example 2:**
100
101
       _Input:_
102
      Triplets from a knowledge graph, about algorithms and their properties.
103
104
      _Output:_
105
         `json
106
      Ε
107
           {
                "question": "What is the worst-case space complexity of merge sort?",
108
                "choices": [
"0(n)",
"0(1)",
109
110
111
                    "O(n log n)",
112
113
                    "O(n^2)"
               ],
"correct": 0,
114
115
                "supporting_triplets": [
116
                     "merge sort--worst-case space complexity->O(n)",
117
                    "merge sort--has space complexity->O(n)",
118
                    "merge sort--space complexity->O(n)",
119
120
                    "quicksort--worst-case space complexity->O(n)",
121
                    "heapsort--worst-case space complexity->O(1)"
               ٦
122
123
           }
      ]
124
125
126
127
      Note how space complexity can be expressed through variations like:
128
      - "worst-case space complexity"
      - "has space complexity"
129
130
      - "space complexity"
131
      # Rules
132
133
134
      1. Only use information explicitly stated in the triplets

    Only reference entities mentioned in the knowledge base
    Treat the knowledge base as the sole source of truth

135
136
137
       4. If a fact isn't explicitly stated in a triplet, don't assume it
138
      5. Consider semantic equivalence of different predicates
139
140
      # Notes
```

- Verify that distractors are plausible to encourage critical thinking. 142

- 143 - Provide triplets from the sources to guide learners in answering the questions. 144
- Return nothing but the json response, formatted to be loaded by Python's json.loads without any modifications. 145

894

896

900

901

902

903

904

905

907

908

909

916

917

918

919

920

921

Where the edges were added to the message with User role in the subject –predicate \rightarrow object format:

Knowledge base: 2 {edges}

141

A.4 Rule-based EQG from Knowledge Graphs

We used the following prompt to rephrase the rule-based MCQs:

Rewrite the following JSON questions and options to make them more readable and user-friendly by:

922

- Simplifying and improving readability. Correcting any grammatically incorrect sentences. 3
- 4 Replacing unnecessarily long date formats with concise and clear ones (e.g.: in the case of Jan 01 dates show only the year).
- Phrasing the question in such a way that it is less like it was generated using a template. 5
- Make the questions more creative, but preserve the original meaning 6

8 Return nothing but the json response, formatted to be loaded by Python's json.loads without any modifications

B **Generated Questions**

Atomic Facts 897

- 1. Which algorithm is a hybrid of merge sort and insertion sort? 925 Heapsort
 - 926 • Timsort 927 Ouicksort
 - 928 • Bubble Sort
 - 2. Which algorithm is most suitable for sorting linked lists?
 - 931 · Merge Sort 932 · Quicksort 933 • Heapsort 934 Selection Sort
- 935 3. Which algorithm is particularly efficient for 910 sorting linked lists? 911 937
- Merge Sort 912 938 • Quicksort 013 939 • Heapsort 914 940 • Bubble Sort 915 941
 - 4. Which sorting algorithm has a worst-case time complexity of $O(n\hat{2})$? 943
 - · Bubble Sort 944 Merge Sort 945 • Heapsort 946 • Timsort 947

- 5. Which sorting algorithm is known for its efficiency in sorting arrays with many equal elements?
 - Quicksort
 - Merge Sort
 - Heapsort
 - Selection Sort
- 6. Which sorting algorithm is known for its poor locality of reference?
 - Heapsort
 - Quicksort
 - Merge Sort
 - Insertion Sort
- 7. Which sorting algorithm is known for its simple implementation but poor performance on large datasets?
 - Merge Sort
 - Quicksort
 - Bubble Sort
 - Timsort
- 8. Which sorting algorithm is known to be stable?
 - Merge Sort
 - Quicksort
 - Heapsort
 - Selection Sort

9. Which sorting algorithm is known to perform poorly on already sorted data?		949 <u>Insertion Sort</u> 949 Quicksort
 Quicksort with the first element as the pivot Merge Sort 		950 [•] Heapsort 951 [•] Bubble Sort
HeapsortInsertion Sort	16. V	Which sorting algorithed algorithed field that the source of the source
10. Which sorting algorithm is known to perform poorly on large lists due to its quadratic time complexity?A Marga Sort		 954 <u>Insertion Sort</u> 955 Selection Sort 956 Merge Sort 956 Ouicksort
 Merge Soft Quicksort <u>Bubble Sort</u> Heapsort 	17. Y	Which sorting algorith alfback when quickso • Heapsort
11. Which sorting algorithm is particularly ben- eficial when sorting data stored on slow-to- access media?		 96• Merge Sort 96• Bubble Sort 96• Insertion Sort
 Quicksort Bubble Sort Merge Sort Insertion Sort 	18. V	which sorting algorith n Python since versio • Quicksort • Heapsort
12. Which sorting algorithm is particularly inefficient for large lists?		969 <u>Timsort</u> 969 Merge Sort
 <u>Bubble Sort</u> Merge Sort Quicksort Heapsort 	19. v	Which sorting algorit conquer' approach? 972 Quicksort 973
13. What is the main advantage of Timsort over Quicksort for sorting object references or pointers?	20.	 Insertion Sort 974 Bubble Sort 975 Selection Sort 976 976 What is the primary a
 Less memory usage Better locality of reference <u>Stability</u> Faster execution time 	S	 soft over Quicksort? 978 978 979 Lower space con Faster average-ca
14. Which sorting algorithm is specifically de- signed to handle large datasets stored on slow- to-access media?		98• Better worst-case 982 983
 Merge Sort Quicksort 		984 985
HeapsortInsertion Sort		986 987
15. Which sorting algorithm is typically less effi- cient on large lists compared to merge sort?		988 989

948	Insertion Sort	990
949	Quicksort	991
95 0	Heapsort	992
951	Bubble Sort	993
. Whi	ch sorting algorithm is typically more ef-	994
ficie	nt than bubble sort on average?	995
954	Insertion Sort	996
955	Selection Sort	997
956	Merge Sort	998
957	Quicksort	999
. Whi fall	ch sorting algorithm is typically used as a	1000
960	uek when queksort becomes memorent.	1001
•	Heapsort	1002
96†	Merge Sort	1003
962	Bubble Sort	1004
	Insertion Sort	1005
. Whi	ch sorting algorithm is used as the default	1006
inP	ython since version 2.3?	1007
966	Quicksort	1008
967	Heapsort	1009
968	Timsort	1010
969	Merge Sort	1011
. Whi	ch sorting algorithm uses a 'divide-and-	1012
conc	juer' approach?	1013
972 •	Quicksort	1014
973	Insertion Sort	1015
974	Bubble Sort	1016
975	Selection Sort	1017
976 . Wha	t is the primary advantage of using Tim-	1018
sort	over Quicksort?	1019
978 •	Stability	1020
979 •	Lower space complexity	1021
980 •	Faster average-case performance	1022
98†	Better worst-case time complexity	1023
982		
983		
984		
985		
986		
987		
988		

1024	KG + LLM 1065
1025	1. Which algorithm has a best-case time com-
1026	plexity of $O(n)$? 1067
1027	• Insertion Sort 1068
1028	Merge Sort 1069
1029	• Heapsort 1070
1030	• Quicksort 1071
	1072
1031	2. Which algorithm has a worst-case space com- playity of $O(1)^2$ 1073
1032	1074
1033	• Heapsort
1034	Merge Sort
1035	• Quicksort
1036	• Timsort 1078
1037	3. Which algorithm has a worst-case time com-
1038	plexity of $O(n\hat{2})$?
1039	Bubble Sort
1040	• Heapsort
1041	Merge Sort
1042	• Timsort
	1084
1043	4. Which algorithm has a worst-case time com-
1044	plexity of $O(n_2)$? 1086
1045	• <u>Bubble Sort</u> 1087
1046	• Timsort 1088
1047	• Heapsort 1089
1048	Merge Sort 1090
1049	5. Which algorithm has an average time com-
1050	plexity of $O(n \log n)$? 1092
1051	• Heapsort
1052	• Bubble Sort
1053	• Insertion Sort 1095
1054	• Selection Sort 1096
	1097
1055	6. Which algorithm is a derivative work of 1im-
1050	1099
1057	• Insertion Sort 1100
1058	• Merge Sort
1059	• Heapsort 1102
1060	• Bubble Sort
1061	7. Which algorithm is an example of an adaptive
1062	sort? 1105
1063	• Timsort 1106
1064	• Heapsort 1107

- Quicksort
- Selection Sort
- 8. Which algorithm is based on both merge sort and insertion sort?
 - Timsort
 - Heapsort
 - Quicksort
 - Bubble Sort
- 9. Which algorithm is based on both merge sort and insertion sort?
 - Timsort
 - Heapsort
 - Quicksort
 - Bubble Sort
- 10. Which algorithm is named after Tim Peters?
 - Timsort
 - Bubble Sort
 - Insertion Sort
 - Merge Sort
- 11. Which algorithm is not a comparison sort?
 - Heapsort
 - Timsort
 - Bubble Sort
 - Merge Sort
- 12. Which algorithm is not a stable sorting algorithm?
 - Heapsort
 - Bubble Sort
 - Insertion Sort
 - Merge Sort
- 13. Which algorithm is not a stable sorting algorithm?
 - Heapsort
 - Timsort
 - Merge Sort
 - Bubble Sort
- 14. Which algorithm is used by Python for sorting?
 - Timsort
 - Quicksort
 - Heapsort
 - Merge Sort

15. Which algorithm is used by Python?

	 <u>Timsort</u> Bubble Sort 	1108
	Insertion Sort	1110
	Heapsort	1111
16.	Which algorithm is used by the Java Platform, Standard Edition?	1112 1113
	• Timsort	1114
	Bubble Sort	1115
	• Insertion Sort	1116
	• Heapsort	1117
17.	Which sorting algorithm has a best-case time complexity of O(n)?	1118 1119
	• Insertion Sort	1120
	• Heapsort	1121
	• Merge Sort	1122
	• Quicksort	1123
18.	Which sorting algorithm has a worst-case	1124
	space complexity of O(1)?	1125
	• Heapsort	1126
	Merge Sort	1127
	• Timsort	1128
	• Quicksort	1129
19.	Which sorting algorithm is an instance of a	1130
	comparison sort?	1131
	• Quicksort	1132
	• Heapsort	1133
	Merge Sort	1134
	• Bubble Sort	1135
20.	Which sorting algorithm is an instance of an	1136
	adaptive sort?	1137
	• <u>Timsort</u>	1138
	• Heapsort	1139
	• Quicksort	1140
	• Merge Sort	1141

1142	KG rules 1182
1143	1. What kind of algorithm is Merge Sort?
1144	Sorting Algorithm
1145	Stable Sorting Algorithm
1146	Online Algorithm ¹¹⁸⁶
1147	• Divide-and-Conquer Algorith ¹¹⁸⁷
1148	2. What kind of algorithm is Selection Sort?
1149	• Sorting Algorithm 1190
1150	 Divide-and-Conquer Algorith¹⁰⁹¹
1151	Online Algorithm 1192
1152	• Adaptive Sort 1193
1153	3. What kind of algorithm is Timsort? ¹⁹⁴
1154	<u>Comparison Sort</u> 1196
1155	• Divide-and-Conquer Algorithm
1156	Online Algorithm
1157	• In-Place Algorithm 1198
1158	4. What type of algorithm is Timsort?
1159	• Adaptive Sort 1201
1160	 Divide-and-Conquer Algorithmo2
1161	• In-Place Algorithm
1162	Online Algorithm 1204
1163	5. What type of algorithm is Timsort?205
1164	• Online Algorithm 1206
1165	• Divide-and-Conquer Algorithm
1166	• In-Place Algorithm
1167	Stable Sorting Algorithm 1209
1168	6. Which algorithm has a best-case time com-
1169	plexity of O(n log n)?
1170	Merge Sort
1171	Bubble Sort
1172	• Heapsort
1173	• Selection Sort
	1216
1174	7. Which algorithm has a worst-case space com- playity of $O(1)$ auxiliary?
1175	plexity of O(1) auxiliary? 1218
1176	• <u>Bubble Sort</u> 1219
1177	• Insertion Sort 1220
1178	• Heapsort
1179	• Timsort 1222
1180	8. Which algorithm has a worst-case space com-
1181	plexity of O(1)?

- Heapsort
- Merge Sort
- Quicksort
- Bubble Sort
- 9. Which algorithm has a worst-case time complexity of O(n log n)?
 - Insertion Sort
 - Quicksort
 - Merge Sort
 - Selection Sort
- 10. Which algorithm has a worst-case time complexity of O(n log n)?
 - Selection Sort
 - Insertion Sort
 - <u>Timsort</u>
 - Bubble Sort
- 11. Which algorithm is a derivative of Timsort?
 - Insertion Sort
 - Selection Sort
 - Quicksort
 - Timsort
- 12. Which algorithm is an example of a divideand-conquer algorithm?
 - Selection Sort
 - Quicksort
 - Heapsort
 - Insertion Sort
- 13. Which algorithm is an example of an online algorithm?
 - Insertion Sort
 - Timsort
 - Merge Sort
 - Bubble Sort
- 14. Which algorithm is based on Insertion Sort?
 - Bubble Sort
 - Timsort
 - Heapsort
 - Merge Sort
- 15. Which algorithm is based on Merge Sort?
 - <u>Timsort</u>
 - Selection Sort
 - Quicksort
 - Insertion Sort

16. Which algorithm is derived from Smoothsort?	1225
Selection Sort	1226
Insertion Sort	1227
Merge Sort	1228
• Heapsort	1229
17. Which algorithm is used by the V8 engine?	1230
Insertion Sort	1231
Bubble Sort	1232
• <u>Timsort</u>	1233
• Heapsort	1234
18. Which algorithm uses the merge algorithm?	1235
Quicksort	1236
<u>Merge Sort</u>	1237
Bubble Sort	1238
Insertion Sort	1239
19. Which sorting method is an example of an in-place algorithm?	1240 1241
Quicksort	1242
Selection Sort	1243
<u>Insertion Sort</u>	1244
• Heapsort	1245
20. Which sorting method is named after bubbles?	1246
• <u>Bubble Sort</u>	1247
Insertion Sort	1248
Selection Sort	1249
• Timsort	1250

1251	Plain text	1293
1252	1. In what scenario is Merge Sor	t particularly
1253	advantageous?	1295
1254	• Sorting small datasets	1296
1255	Sorting large datasets with	external storage
1256	 Sorting data in-place 	1298
1257	• Sorting data with minimal	memory 1300
1258	2. What is a common optimization	for auicksort
1259	to avoid worst-case performanc	e? 1302
1260	• Using a fixed pivot	1303
1261	• Using insertion sort for sm	all arrays
1262	• Switching to bubble sort	1305
1263	 Increasing recursion depth 	1306
1264	3. What is a key advantage of	Timsort over
1265	Quicksort?	1308
1266	• Timsort is a stable sort	1309
1267	• Timsort has better worst-c	ase time com-
1268	plexity	1011
1269	• Timsort uses less memory	1312
1270	• Timsort is easier to implen	nent313
1271 1272	4. Which algorithm is described as element to partition the array int	using a 'pivot' to sub-arrays?
1273	Insertion Sort	1316
1274	Merge Sort	1317
1275	Ouicksort	1318
1276	Selection Sort	1319
1277	5. What is a key disadvantage of I	1320 neapsort com- 1321
1210		1322
1279	• Heapsort is not stable.	1323
1280	• Heapsort has a higher wo	1324 time
1000	 Heapsort is not an in-place 	algorithm
1283	Heapsort is recursive	1326
1205	ricapsort is recursive.	1327
1284	6. What is a primary disadvantage	e of Selection
1285	Sort compared to Insertion Sort	? 1329
1286	• Higher time complexity	1330
1287	• More memory usage	1331
1288	• Less efficient for small dat	aset _{§32}
1289	• More write operations	1333
1290	7. Which sorting algorithm is desc	ribed as a sta-
1291	ble, hybrid sorting algorithm t	hat combines
1292	merge sort and insertion sort?	1336

•	Heapsort
---	----------

- <u>Timsort</u>
- Quicksort
- Bubble Sort
- 8. Which sorting algorithm is described as having a worst-case time complexity of $O(n\hat{2})$ but is simple and can be advantageous when auxiliary memory is limited?
 - Selection Sort
 - Merge Sort
 - Heapsort
 - Timsort
- 9. Which sorting algorithm is generally more efficient for partially sorted data?
 - Heapsort
 - Bubble Sort
 - Insertion Sort
 - Selection Sort
- 10. Which sorting algorithm is known for its efficient performance on small datasets and partially sorted data?
 - Insertion Sort
 - Bubble Sort
 - Merge Sort
 - Quicksort
- 11. Which sorting algorithm is known for its poor locality of reference?
 - Heapsort
 - Merge Sort
 - Quicksort
 - Timsort
- 12. What is a significant disadvantage of heapsort compared to quicksort?
 - It is not stable
 - It has a higher worst-case time complexity
 - It uses more memory
 - It is harder to implement
- 13. Which sorting algorithm is known for its stable sorting and efficient performance on sequentially accessed data?
 - Heapsort
 - Quicksort
 - Merge Sort

- Selection Sort
- 14. Which sorting algorithm is noted for its inefficiency on large datasets but is simple and often used for educational purposes?
 - Merge Sort
 - Bubble Sort
 - Quicksort
 - Insertion Sort
- 15. Which sorting algorithm is particularly inefficient for large datasets due to its $(O(n^2))$ complexity?
 - Bubble Sort
 - Heapsort
 - Timsort
 - Merge Sort
- 16. Which sorting algorithm is particularly noted for its use in educational contexts due to its simplicity, despite its inefficiency?
 - Heapsort
 - Selection Sort
 - Bubble Sort
 - Insertion Sort
- 17. Which sorting algorithm is typically faster for randomized data due to better cache coherence?

	133• Heapsort	1362
	• Merge Sort	1363
	1339 Quicksort	1364
	1340 Bubble Sort	1365
18.	Which sorting algorithm is typically used in	1366
	hybrid forms like Timsort due to its stability	1367
	and efficiency for sequential media?	1368
	¹³⁴⁴ Merge Sort	1369
	1349 Quicksort	1370
	134• Heapsort	1371
	¹³⁴⁷ Selection Sort	1372
19.	Which sorting algorithm is used in Python due	1373
	to its efficiency on real-world data?	1374
	135 [•] Heapsort	1375
	• Merge Sort	1376
	135 ^o <u>Timsort</u>	1377
	Bubble Sort	1378
20.	What is the primary advantage of merge sort	1379
	over quicksort?	1380
	¹³⁵⁷ Merge sort is in-place.	1381
	¹³⁵⁸ Merge sort is stable.	1382
	1359 Merge sort has better average time com-	1383
	1360 plexity.	1384
	136• Merge sort requires less space.	1385