

# JudgeAgent: Beyond Static Benchmarks for Knowledge-Driven and Dynamic LLM Evaluation

Anonymous ACL submission

## Abstract

Current evaluation methods for large language models (LLMs) primarily rely on static benchmarks, presenting two major challenges: limited knowledge coverage and fixed difficulties that mismatch with the evaluated LLMs. These limitations lead to superficial assessments of LLM knowledge, thereby impeding the targeted model optimizations. To bridge this gap, we propose JudgeAgent, a knowledge-driven and dynamic evaluation framework for LLMs. To address the challenge of limited knowledge coverage, JudgeAgent leverages LLM agents equipped with context graphs to traverse knowledge structures systematically for question generation. Furthermore, to mitigate data contamination and difficulty mismatch, it adopts a difficulty-adaptive and multi-turn interview mechanism. Thereby, JudgeAgent can achieve comprehensive evaluations and facilitate more effective improvement of LLMs. Empirical results demonstrate that JudgeAgent enables more comprehensive evaluations and facilitates effective model iterations, highlighting the potential of this knowledge-driven and dynamic evaluation paradigm. The source code is available on <https://anonymous.4open.science/r/JudgeAgent>.

## 1 Introduction

Evaluating Large Language Models (LLMs) to verify their knowledge is a crucial step for their successful application across various domains.(Tang et al., 2024; Yuan et al., 2023; Shi, 2024; Wang et al., 2025). Current evaluation methods mainly rely on static benchmarks(Clark et al., 2018; Hendrycks et al., 2021; Huang et al., 2023; Lin et al., 2022; Cobbe et al., 2021; Tang and Yang, 2024), where LLMs are evaluated by their performance on answering predefined static questions. In the early stages, benefiting from the controllable question quality and rapid workflows of static benchmarking, developers can quickly grasp the

core competencies of evolving LLMs, thereby accelerating iteration cycles. However, due to the static nature(Gu et al., 2024; Ko et al., 2024), these benchmarks have become saturated more and more quickly in recent years. For example, it took 3 years for MMLU(Hendrycks et al., 2020) to reach 80% accuracy record, while GPQA(Rein et al., 2023) achieved the same record just within one year. The static nature restricts the evaluations within a pre-defined and limited knowledge scope(Wang et al., 2023a; Kwan et al., 2024), and increases the risk of data contamination(Schaeffer, 2023; Oren et al., 2023), resulting in LLMs being misapplied in scenarios out of their knowledge.

Consequently, researchers attempted to modify static benchmarks using LLMs(Bai et al., 2023; Shi et al., 2025). However, they still face two major challenges. **(1) Limited knowledge coverage:** Constrained by the evaluator LLM’s knowledge limitations and inherent randomness, the knowledge extended by the LLM from static benchmarks remains uncontrollable and limited, making it difficult to achieve comprehensive knowledge evaluations. **(2) Mismatch between question difficulty and the evaluated LLMs:** Current mainstream methods lack adaptive adjustments to question difficulty. Thus, the question difficulty may significantly deviate from the actual knowledge mastery of the evaluated LLMs, making the evaluations harder to precisely reflect the model’s capabilities. Thus, a new knowledge-driven and dynamic evaluation framework is needed.

In this paper, we propose JudgeAgent, a knowledge-driven and dynamic evaluation framework for LLMs. To address the challenge of limited knowledge coverage, JudgeAgent leverages LLM agents with context graphs to traverse knowledge structures at greater breadth and depth. Through graph sampling strategies and knowledge-driven data synthesis, JudgeAgent generates questions with broader and deeper knowledge cover-

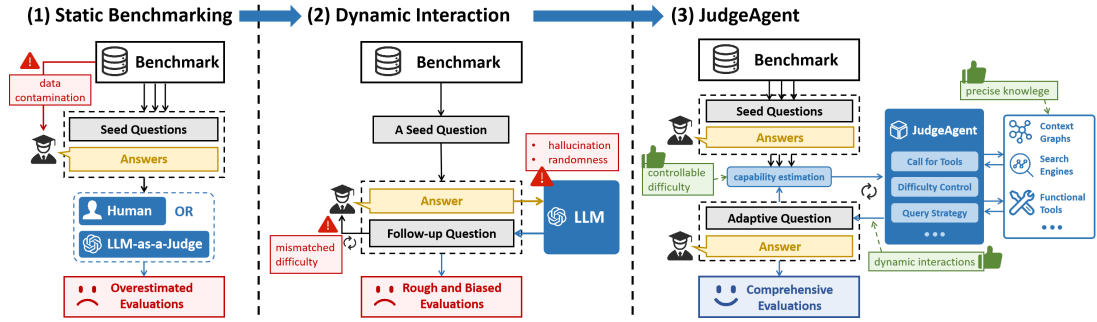


Figure 1: The difference between JudgeAgent and current evaluation paradigms.

age, further exploring potential knowledge deficiencies in evaluated LLMs. Furthermore, to mitigate data contamination and difficulty mismatch, JudgeAgent dynamically adjusts the capability estimation and generates difficulty-adaptive questions based on the target’s responses in multi-turn interviews. Thereby, JudgeAgent can provide comprehensive evaluations and facilitate more effective subsequent improvement of evaluated LLMs.

Extensive experiments and analysis validate that JudgeAgent can provide effective and precise evaluation results. Our in-depth analysis also reveals the significant potential of this evaluation paradigm. It offers valuable insights into how additional knowledge-driven and dynamic strategies for evaluator LLMs can be designed, which enhance the quality of dynamic evaluation results.

In summary, our contributions are as follows:

- We introduce JudgeAgent, a knowledge-driven and dynamic evaluation framework for LLMs. JudgeAgent leverages LLM agents with context graphs to traverse knowledge structures and perform dynamic multi-turn interviews, achieving comprehensive knowledge evaluations.
- We propose a difficulty-adaptive and multi-turn evaluation mechanism that simulates an interactive interview, enabling the evaluations to more precisely reflect the knowledge mastery of the evaluated LLM.
- We conduct extensive experiments to validate the effectiveness of JudgeAgent and highlight its significant potential for the dynamic evaluation paradigm.

## 2 Related Works

### 2.1 Static Benchmark-based Evaluation

These methods employ pre-constructed benchmarks to evaluate LLMs, using formats such

as multiple-choice, question-answer(Q&A), or prompts for performing tasks. For multiple-choice (Clark et al., 2018; Hendrycks et al., 2021; Huang et al., 2023) or Q&A (Lin et al., 2022; Cobbe et al., 2021; Tang and Yang, 2024), the benchmark provides correct answers and assesses the target by the accuracy. For task-execution format, the benchmark measures the model by the metrics assessed by human (Chang et al., 2024) or LLM-as-a-judge (Wang et al., 2023b; Zheng et al., 2023). These methods ensure controllable question quality, but they are also susceptible to data contamination, which undermines the validity of evaluations.

### 2.2 Dynamic LLM-based Evaluation

With the deepening application of LLMs, in addition to mainstream static methods, recent researchers have begun to leverage LLMs to dynamically generate evaluation questions. Bai et al. (2023) employs LLMs as examiners to generate questions based on Google Trends categories. SafetyQuizzer(Shi et al., 2025) formulates questions with current events from search engines to maximize the timeliness. KIEval(Yu et al., 2024), and LLM-as-an-Interviewer (Kim et al., 2025) leverage LLMs to generate follow-up questions based on the target’s responses in multi-turn interactions.

Compared to static methods, these dynamic methods address the challenge of data contamination and provide more genuine evaluations. However, the generations of these methods rely mainly on the LLM, resulting in limited knowledge coverage and pool quality control. Moreover, these methods lack adaptive adjustments of question difficulty, resulting in a deviation between the evaluations and the actual capabilities of target LLMs. This mismatch leads to biased subsequent optimization efforts and reduced efficiency.

Thus, we propose a knowledge-driven and dynamic evaluation framework named JudgeAgent. To expand knowledge coverage, JudgeAgent lever-

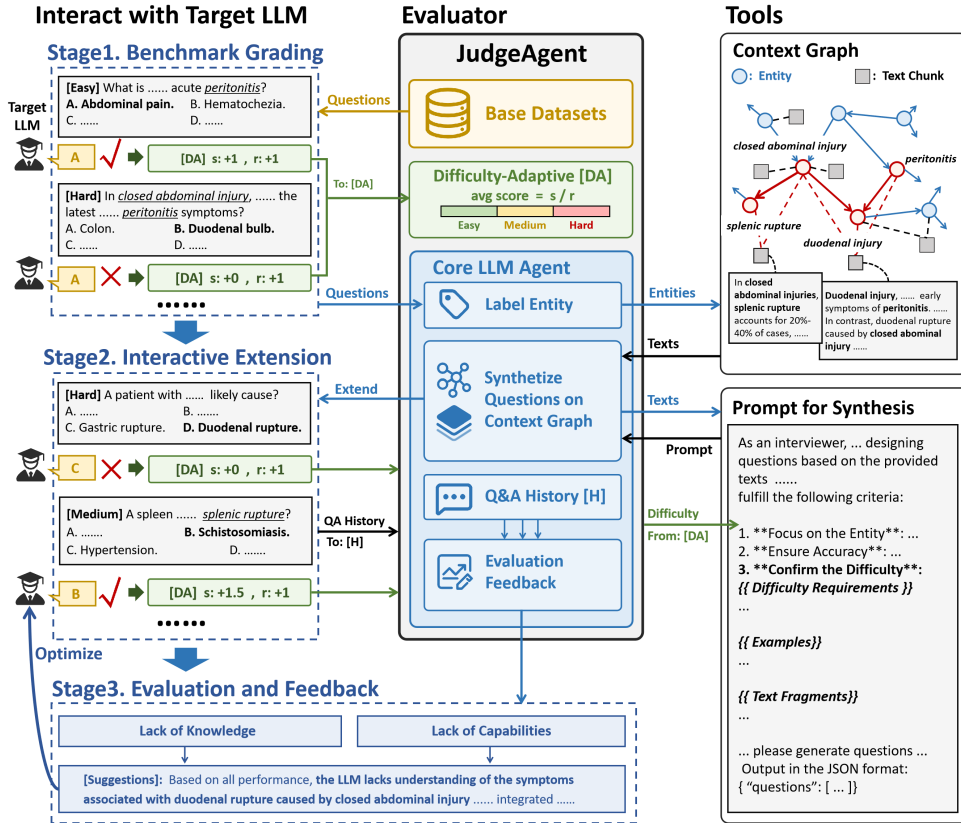


Figure 2: The framework of JudgeAgent. The left part is the interaction process. The central part is the composition of JudgeAgent. The right part presents the tools of JudgeAgent.

ages LLM agents with context graphs to traverse knowledge structures and generate questions. To mitigate data contamination and difficulty mismatch, JudgeAgent adopts a difficulty-adaptive multi-turn interview mechanism to conduct dynamic evaluations. Thus, JudgeAgent provides comprehensive evaluations and facilitates subsequent improvement of LLMs.

### 3 Methodology

We introduce JudgeAgent, a knowledge-driven dynamic evaluation framework, which utilizes LLM agents with context graphs to simulate the entire interview of target LLMs. The workflow comprises three core components: (1) **Benchmark Grading**: Fundamentally assessing target LLMs through testing on public static benchmarks. (2) **Interactive Extension**: Using knowledge-driven synthesis to dynamically generate follow-up questions and update the difficulty based on the estimations of target LLMs. (3) **Evaluation Feedback**: Evaluating the target’s deficiencies in multiple dimensions and providing actionable suggestions. The overall framework is illustrated in Figure 2, and detailed as pseudo code in Algorithm 1.

#### 3.1 Benchmark Grading

To address the challenge of difficulty mismatch in prior dynamic evaluation, JudgeAgent first evaluates the target LLM on public static benchmarks and obtains a base score through a difficulty scoring mechanism. Similar to a written test before an interview, this allows JudgeAgent to form a basic capability estimation of the target LLM.

JudgeAgent first partitions the benchmark questions into batches. Then it prompts the target LLM to answer the questions batch by batch. After assessing the responses, JudgeAgent estimates the target’s capability based on its performance in each batch. Inspired by the common real-world practice of scoring student exams, we design a linear difficulty control mechanism. Within each batch, each correct answer improves the total score, and the average score determines the follow-up question’s difficulty. Let  $[d_1, d_2, d_3]$  respectively denote the difficulty *Easy*, *Medium*, and *Hard*, for the score gain  $c_{it}$  for a correct answer at different difficulty levels  $d_i$ , we set  $[c_{1t}, c_{2t}, c_{3t}] = [1, 1.5, 2]$ . And for the threshold  $t_{ij}$  between  $d_i$  and  $d_j$ , we set  $[t_{12}, t_{23}] = [0.5, 1]$ . Once the dynamic average score exceeds the threshold  $t_{ij}$ , the difficulty of

211	the next generated question will be set to $d_j$ . The	edge mastery. For instance, when an overly chal-	262
212	mechanism’s detailed design, analysis, and proofs	lenging test evaluates a model with limited knowl-	263
213	are available in Appendix B.	edge, the evaluation may reveal that the model per-	264
214	<b>3.2 Interactive Extension</b>	forms poorly, but fails to delineate the actual range	265
215	This process is an "interview" where JudgeAgent	of its knowledge. Therefore, JudgeAgent dynam-	266
216	conducts a knowledge-driven and difficulty-	ically adjusts the follow-up question’s difficulty	267
217	adaptive multi-turn evaluation. Specifically,	based on capability estimation through the diffi-	268
218	JudgeAgent dynamically adjusts the difficulty and	culty control mechanism in Section 3.1.	269
219	generates questions with expanded knowledge	Based on the capability estimation, JudgeAgent	270
220	based on the responses. Furthermore, JudgeAgent	adjusts the question difficulty as <i>Easy</i> , <i>Medium</i> ,	271
221	explores the potential knowledge gaps of target	or <i>Hard</i> , and explicitly specifies each difficulty’s	272
222	LLMs through a multi-turn interaction process.	requirements in the prompt. For <i>Easy</i> , the focus is	273
223	Each iteration in this process consists of three steps:	on assessing information retrieval skills, primarily	274
224	Relevant Knowledge Retrieval, Difficulty-Adaptive	through cloze-style questions based on the original	275
225	Question Generation, and Capability Estimation.	text, like " <i>Which virus primarily causes HFMD?</i> ".	276
226	<b>Relevant Knowledge Retrieval:</b> To achieve	For <i>Medium</i> , the emphasis is on understanding key	277
227	knowledge-driven evaluation, JudgeAgent invokes	concepts, with questions involving basic inference,	278
228	external knowledge and constructs context graphs	like " <i>Which description is correct regarding meio-</i>	279
229	to facilitate the association and extension of knowl-	<i>sis II?</i> " For <i>Hard</i> , the questions are designed to	280
230	edge during evaluation. If only reference texts of	encourage deep thinking and complex logical anal-	281
231	seed questions are used for generation, the gener-	ysis, like " <i>A patient ingested a toxic substance...</i>	282
232	ated questions will be highly similar to the origi-	<i>The doctor employed... which type of treatment is</i>	283
233	nal ones. Inspired by SoG(Jiang et al., 2025),	<i>likely lacking?</i> ". Detailed prompt designs are pro-	284
234	JudgeAgent employs a context-graph-based sam-	vided in Appendix G, and the validation of question	285
235	pling approach to get background texts. The con-	quality is described and analyzed in Appendix E.4.	286
236	text graph is constructed from the benchmark’s	<b>Capability Estimation:</b> JudgeAgent queries the	287
237	knowledge base with entities and chunks extracted	target LLM with generated questions. Based on	288
238	from texts as nodes. The entities in the same text	LLM’s responses, JudgeAgent computes average	289
239	will be linked on the graph. The detailed construc-	difficulty scores with the strategy in Section 3.1,	290
240	tion is illustrated in Appendix C and Algorithm 3.	and determines the difficulty for the next round’s	291
241	During sampling, JudgeAgent first extracts enti-	generation. Through this adaptive multi-turn mech-	292
242	tities from seed questions and finds the most similar	anism, JudgeAgent progressively aligns the assess-	293
243	entity $e$ on the graph. In each hop, JudgeAgent	ment with the true capabilities of the target LLM,	294
244	retrieves the most relevant chunks of $e$ to the ques-	thereby achieving more precise evaluations.	295
245	tion, and randomly selects a chunk $c$ . $e$ is the seed	<b>3.3 Evaluation Feedback</b>	296
246	entity of the next hop, and $c$ is added to the path.	Inspired by Generative Reward Model methods	297
247	After $N$ hops, we sample a path with breadth and	(Ankner et al., 2024; Ye et al., 2024; Li et al., 2025)	298
248	depth, and the chunks on the path are concatenated	and the evaluation report method in (Kim et al.,	299
249	to form the knowledge for generation. The context	2025), JudgeAgent generates a text evaluation re-	300
250	graph ensures both knowledge breadth and depth,	port based on all Q&A history for each batch. Ben-	301
251	and the greedy similarity strategy ensures relevance	efiting from context graphs, the evaluation report	302
252	to seed questions. Entity extraction is performed	can identify the knowledge deficiencies of target	303
253	by GPT-4.1, with prompts detailed in Appendix G.	LLMs and provide concise and valuable feedback	304
254	<b>Difficulty-Adaptive Question Generation:</b> To	around key knowledge concepts.	305
255	mitigate data contamination and difficulty mis-	<b>4 Experiments</b>	306
256	match, JudgeAgent conducts multi-turn interac-	In this section, we validate the effectiveness of	307
257	tive interviews to explore potential deficiencies and	JudgeAgent using the method in Section 3.3, which	308
258	incorporates adaptive difficulty adjustments into	is detailed in Algorithm 2. The following re-	309
259	question generation. A fixed-difficulty test can	search questions guide our experiments: <b>RQ1</b> ,	310
260	identify underperforming models, but can hardly		
261	comprehensively evaluate a model’s actual knowl-		

311 does JudgeAgent’s knowledge-driven evaluation  
312 genuinely discover the shortcomings of the target  
313 model? **RQ2**, is JudgeAgent’s dynamic evaluation  
314 more effective and precise than the static bench-  
315 marking? **RQ3**, to what extent does each mecha-  
316 nism in JudgeAgent influence the evaluations?

317 Current LLM knowledge evaluation methods  
318 lack quantifiable and standardized validation ap-  
319 proaches. Therefore, to validate the effectiveness  
320 of JudgeAgent, we use a similar approach as (Kim  
321 et al., 2025). We use the evaluation report produced  
322 by JudgeAgent as suggestions to prompt the target  
323 LLM to answer the same questions again. Then  
324 we compare the accuracy before and after the inter-  
325 vention to indirectly validate the effectiveness. To  
326 avoid cheating, we produce suggestions without di-  
327 rectly providing the correct answer or background  
328 knowledge of seed questions to JudgeAgent.

329 **Experiment Setup.** We select GPT-4.1<sup>1</sup> to be  
330 the core LLM. In our experiments, the batch size  
331 in Benchmark Grading is 3, and the Interactive  
332 Extension is limited to a maximum of 3 rounds.

#### 333 4.1 Dataset and Target Model Selection

334 In our experiments, we select MedQA(Jin et al.,  
335 2021), MultiHop-RAG(Tang and Yang, 2024), and  
336 QuALITY(Pang et al., 2022) as the initial bench-  
337 marks. We remove the background knowledge  
338 of questions from MedQA and MultiHop-RAG  
339 to evaluate the target’s knowledge rather than  
340 its comprehension ability, and to verify whether  
341 JudgeAgent can discover the target’s actual knowl-  
342 edge deficiencies. We use QuALITY to validate  
343 the effectiveness in guiding comprehension and  
344 reasoning. The detailed information and the pre-  
345 processing procedures are in Appendix D.

346 For target LLMs, we select Qwen3(Yang  
347 et al., 2025), GLM4-Flash(GLM et al., 2024),  
348 GPT-4.1(Achiam et al., 2023), and gemini-2.5-  
349 pro(Comanici et al., 2025)<sup>2</sup>. We utilize the official  
350 APIs to interact with these models.

#### 351 4.2 Evaluation Metrics

352 The metrics used in our experiments are as follows:  
353 (1) **Accuracy (ACC).** We use this metric to mea-  
354 sure the performance of the target model on  
355 base datasets. To investigate the effectiveness of

JudgeAgent, we compared the changes in the ACC  
before and after dynamic evaluation. ACC1 repre-  
sents the target’s ACC in answering base questions  
in *Benchmark Grading*, measuring its performance  
before receiving evaluation feedback. ACC2 de-  
notes the ACC in answering the same questions  
after receiving feedback from JudgeAgent.

(2) **Correction Rate (CR).** This metric quanti-  
fies the proportion of questions that the target  
model initially answered incorrectly but subse-  
quently answered correctly after receiving evalua-  
tion suggestions, which measures the effectiveness  
of JudgeAgent. A higher correction rate indicates  
better performance of the suggestions.

(3) **Correct-to-Error Rate (CtE).** This metric  
serves as the inverse of the Correction Rate. A  
lower CtE indicates greater effectiveness.

#### 373 4.3 Main Results and Analysis

374 To address **RQ1**, we conduct experiments on  
375 three datasets to validate the effectiveness of  
376 JudgeAgent’s evaluation, and the results are shown  
377 in Table 1 and Table 2.

378 Based on the results from MedQA and Multi-  
379 HopRAG in Table 1, JudgeAgent can effectively  
380 identify potential knowledge deficiencies in target  
381 LLMs and subsequently mitigate these deficien-  
382 cies by providing targeted suggestions to the target  
383 LLM. Furthermore, as evidenced by the overall  
384 performance on QuALITY in Table 2, JudgeAgent  
385 also contributes to providing further guidance to ad-  
386 dress potential shortcomings in the logical reason-  
387 ing abilities of the target LLM, thereby assisting in  
388 refining the target’s thinking steps. Additionally, by  
389 comparing CR and CtE of different LLMs before  
390 and after receiving suggestions, it can be observed  
391 that the effectiveness of suggestions is less con-  
392 sistent for relatively weaker LLMs (e.g., the free  
393 model GLM4-Flash), as reflected in the higher CtE.  
394 In contrast, stronger LLMs are less susceptible to  
395 misleading suggestions.

396 Based on the performance across different dif-  
397 ficulties of QuALITY in Table 2, we can ana-  
398 lyze the effectiveness of JudgeAgent for different  
399 targets on various difficulty levels. For stronger  
400 LLMs (Qwen3, GPT-4.1, and Gemini-2.5-pro),  
401 JudgeAgent provides stronger guidance, particu-  
402 larly on questions of greater difficulty (*Medium*  
403 and *Hard*), since these LLMs may have already  
404 mastered the basic knowledge of *Easy* questions.  
405 In contrast, for weaker LLMs, JudgeAgent leads to  
406 considerable improvement across all difficulty lev-

<sup>1</sup>We use gpt-4.1-2025-04-14 from OpenAI’s official API for all experiments

<sup>2</sup>We use qwen-plus-2025-04-28, gpt-4.1-2025-04-14, gemini-2.5-pro-preview-06-05, and the free version glm-4-flash-250414 as the target models.

Target Model	MedQA				MultiHopRAG			
	ACC1	ACC2	CR↑	CtE↓	ACC1	ACC2	CR↑	CtE↓
Qwen3	91.71	96.38	5.02	0.35	63.65	70.07	17.49	11.07
GLM4-Flash	80.09	92.82	13.46	0.73	51.25	65.92	24.26	9.59
GPT-4.1	84.97	92.44	7.65	0.18	68.94	75.55	12.36	5.75
Gemini-2.5-pro	91.04	94.60	4.03	0.47	62.25	71.83	15.41	5.83

Table 1: The results on MedQA and MultiHopRAG, and all values are percentages.

Target Model	QuALITY-overall				QuALITY-easy			
	ACC1	ACC2	CR↑	CtE↓	ACC1	ACC2	CR↑	CtE↓
Qwen3	87.83	88.84	1.88	0.87	94.63	95.61	0.98	0.00
GLM4-Flash	73.48	76.38	4.78	1.88	83.26	84.65	4.19	2.79
GPT-4.1	89.13	92.75	3.91	0.29	93.66	96.10	2.44	0.00
Gemini-2.5-pro	93.77	96.38	3.19	0.58	97.55	99.02	1.47	0.00

Target Model	QuALITY-medium				QuALITY-hard			
	ACC1	ACC2	CR↑	CtE↓	ACC1	ACC2	CR↑	CtE↓
Qwen3	88.53	88.89	1.43	1.08	80.10	82.04	3.40	1.46
GLM4-Flash	74.35	79.18	5.58	0.74	62.14	64.08	4.37	2.43
GPT-4.1	91.04	94.27	3.58	0.36	82.04	87.38	5.83	0.49
Gemini-2.5-pro	93.57	95.36	2.50	0.71	90.29	95.15	5.83	0.97

Table 2: The results on QuALITY with different difficulty levels. QuALITY-X refers to the subdataset that consists of questions with specific difficulty, and -overall refers to all questions.

els, with particularly notable gains on *Easy* level, since JudgeAgent’s feedback is more effective at filling gaps in basic concepts. These results further indicate that JudgeAgent more accurately assesses the capability and provides more difficulty-adaptive guidance to the target LLMs. Furthermore, for all target LLMs, JudgeAgent’s suggestions demonstrate clear guidance on *Hard* questions, indicating that JudgeAgent effectively identifies underlying deficiencies in target LLMs through dynamic interactive evaluation.

In summary, JudgeAgent can effectively and precisely identify potential knowledge deficiencies in target LLMs, offering more refined evaluation results. Furthermore, we conducted cross-validation experiments to verify that the JudgeAgent’s evaluations are effective not only for seed questions, which are detailed in Appendix E.5. Additionally, experiments on the efficiency and stability of JudgeAgent are detailed in Appendix E.1 and E.2.

#### 4.4 Analysis of JudgeAgent’s Dynamic Scores

To address RQ2, we analyze the dynamic difficulty scores and the difficulty distribution of generated questions during JudgeAgent’s Interactive Expansion stage. The results are shown in Table 3.

As the results on MedQA show, there is no significant difference in the accuracy and scores between GPT-4.1 and GLM4-Flash on the static benchmark

(*base*). However, through multi-round dynamic evaluation, it can be observed that the difficulty of questions queried to GPT-4.1 continuously increases, while that for GLM4-Flash decreases, indicating that GPT-4.1’s actual knowledge mastery exceeds the preset difficulty of the static benchmark, further suggesting that GPT-4.1’s knowledge mastery is far higher than that of GLM4-Flash.

Similarly, on MultiHopRAG, Qwen3 and Gemini-2.5-pro show comparable performance on the static benchmark. While in dynamic evaluations, Qwen3’s performance is significantly lower than that of GPT-4.1 and Gemini-2.5-pro, with the gap progressively widening over rounds. This gap is also reflected in the fact that the generated questions for Qwen3 are primarily in *Medium* level. Meanwhile, although GPT-4.1 performs noticeably better than Gemini-2.5-pro on the static benchmark, their performance was comparable in dynamic evaluations, as can be seen from their similar difficulty distributions in the questions generated. Moreover, the multi-turn interview mechanism makes the differentiation in evaluation more pronounced. In summary, compared to static benchmarking, JudgeAgent’s dynamic evaluation can more precisely distinguish the actual knowledge mastery of different LLMs. Furthermore, we validate JudgeAgent’s resistance to data contamination through dynamic evaluation, with experiments and

Benchmark	Target Model	ACC(%)	Difficulty Score				Difficulty Distribution(%)		
			base	@1	@2	@3	easy	medium	hard
MedQA	GLM4-Flash	80.09	1.201	1.230	1.217	1.192	20.86	12.64	66.50
	GPT-4.1	84.97	1.275	1.388	1.395	1.401	14.61	9.02	76.38
	Qwen3	91.71	1.376	1.472	1.496	1.487	8.01	8.58	83.42
	Gemini-2.5-pro	91.04	1.366	1.484	1.494	1.505	9.41	7.76	82.84
MultiHopRAG	GLM4-Flash	51.25	0.769	0.587	0.535	0.486	34.06	48.13	17.81
	GPT-4.1	68.94	1.034	1.006	1.000	0.996	11.55	40.16	48.29
	Qwen3	63.65	0.955	0.801	0.767	0.746	19.62	48.91	31.47
	Gemini-2.5-pro	62.25	0.934	0.967	1.005	1.019	13.63	38.80	47.57
QuALITY	GLM4-Flash	74.35	1.065	1.143	1.232	1.304	12.09	39.30	48.60
	GPT-4.1	91.04	1.320	1.428	1.522	1.603	1.86	23.49	74.65
	Qwen3	88.53	1.296	1.388	1.466	1.520	3.26	25.58	71.76
	Gemini-2.5-pro	93.57	1.397	1.502	1.595	1.686	0.00	17.91	82.09

Table 3: The dynamic difficulty scores and difficulty distribution of questions generated by JudgeAgent. *base* means the score for base static benchmark questions, and *@K* means the scores for extended questions at the K-th iteration.

analysis detailed in Appendix E.3.

#### 4.5 Ablation Study

To address **RQ3**, we conduct ablation studies on MedQA with GLM4-Flash as the target LLM. When maintaining the responses to seed questions unchanged, we removed different modules of JudgeAgent. By comparing the results under different ablation settings, we investigate how much different modules influence the evaluations. The results are shown in Figure 3.

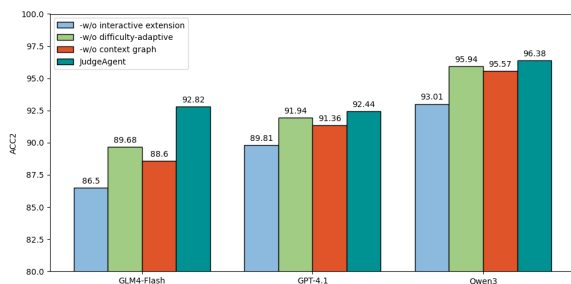


Figure 3: The results of the ablation study with MedQA as the base dataset, and all values are percentages.

JudgeAgent (*w/o* context graph) removes the context graph and only uses chunks sampled randomly from the knowledge base. The results indicate that removing the context graph reduces the effectiveness of the JudgeAgent. The lack of context graph severs the knowledge link between extended questions and base ones, leading to feedback without accurate information, which may disrupt the thinking of target LLMs. The higher CtE compared to the setting *-w/o* difficulty-adaptive also provides supporting evidence for this potential interference.

JudgeAgent (*w/o* difficulty-adaptive) removes the difficulty-adaptive mechanism and generates questions with fixed difficulty rules in the prompt.

The results show that the removal of this module diminishes the effectiveness of JudgeAgent’s evaluations, demonstrating the importance of the difficulty-adaptive mechanism for providing effective evaluations.

JudgeAgent (*w/o* interactive extension) removes the Interactive Extension and only evaluates models with base benchmarks. The results show that removing this stage significantly weakens the effectiveness of JudgeAgent compared to other ablation settings, indicating that the dynamic expansion, which expands both breadth and depth of knowledge, is crucial to the evaluation’s effectiveness.

#### 4.6 Parameter Analysis

The number of extension rounds is crucial for balancing the efficiency and effectiveness of JudgeAgent. We test the JudgeAgent’s effectiveness under different rounds in Interactive Extension, as shown in Figure 5. The results demonstrate that as the rounds increase, the JudgeAgent’s effectiveness, which is represented by the accuracy improvement of target LLMs, also gradually improves. However, the trend gradually slows down, showing a relatively clear marginal effect.

As the rounds increase, JudgeAgent can expand more questions related to the knowledge around base questions. After reaching certain rounds, the Q&A history is sufficient to identify the target’s knowledge deficiencies around base questions. Additional questions serve only as corroborative rather than critical evidence. There is a clear marginal effect in Figure 5. For example, the results of GPT-4.1 on MultiHop-RAG, the improvement per round drops from 2.04% to 1.13%, then to 0.71%, and eventually to 0.23%, and other curves exhibit sim-

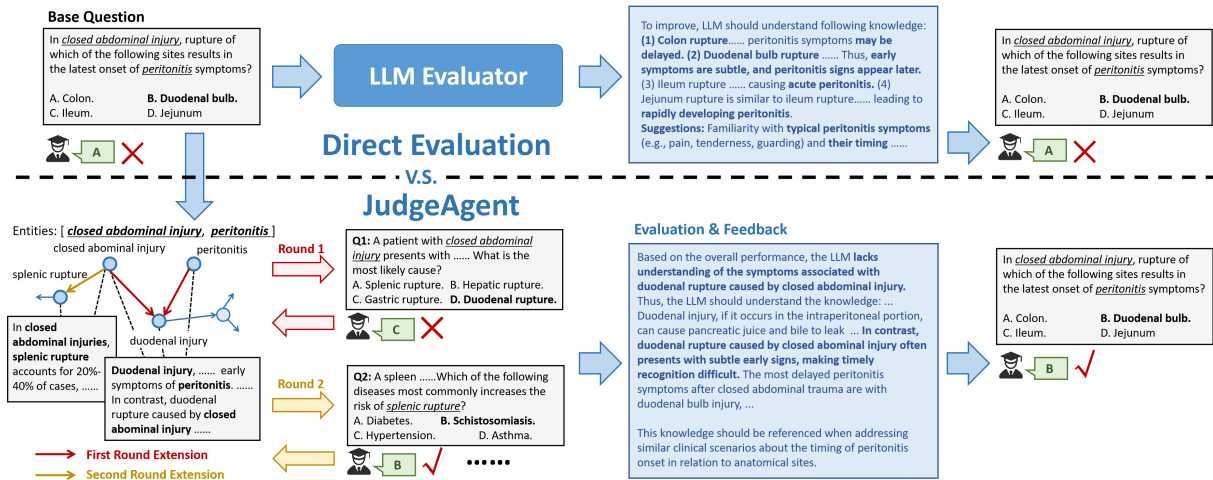


Figure 4: The brief overview of the comparative case in the Case Study.

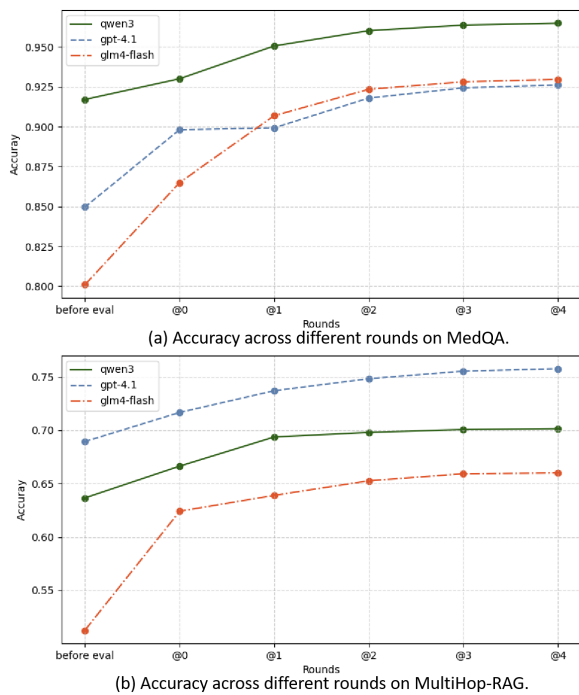


Figure 5: The results of different expansion rounds on MedQA and MultiHop-RAG. @K indicates the ACC improvement after the K-th interaction.

ilar patterns. Therefore, JudgeAgent expands a maximum of 3 rounds in our experiments to avoid wasting resources and time.

Additionally, parameter analysis experiments regarding batch size are presented in Appendix E.6.

#### 4.7 Case Study

To further understand JudgeAgent’s mechanism, we analyze GLM4-Flash’s responses to MedQA questions by receiving suggestions from direct evaluation and JudgeAgent. The seed questions where the target answers correctly are omitted in this case. A brief overview of the case is shown in Figure 4,

and the detailed content is provided in Appendix H.

In this case, the target LLM answers the base question incorrectly. When directly evaluated, the evaluator LLM generically enumerated the potential consequences of closed abdominal injury. Consequently, even after receiving the feedback, the target LLM still outputs an incorrect answer.

In contrast, JudgeAgent first extracted two key entities, *closed abdominal injury* and *peritonitis*. Based on the sampled paths on the context graph, JudgeAgent generated extended questions, which help discover the lack of understanding of *duodenal injury* and *closed abdominal injury*. Ultimately, the target LLM successfully answered the original question correctly with this feedback.

## 5 Conclusion

In this paper, we propose JudgeAgent, a knowledge-driven and dynamic evaluation framework for LLMs. To address the challenge of limited knowledge coverage, JudgeAgent leverages LLM agents with context graphs to traverse knowledge structures and employs knowledge-driven synthesis for question generation. Furthermore, to mitigate data contamination and difficulty mismatch, JudgeAgent dynamically adjusts the capability estimation and generates difficulty-adaptive questions based on the target’s responses in multi-turn interviews. Extensive experiments validate the effectiveness of JudgeAgent and highlight the significant potential of this knowledge-driven and dynamic evaluation paradigm in enhancing the evaluation results. In our future work, we will further refine this dynamic evaluation paradigm and develop more effective and reliable evaluation tools.

## 569 Limitations

570 In this paper, we validate JudgeAgent’s effective-  
571 ness and its great potential for enhancing the dy-  
572 namic evaluation paradigm of LLMs through ex-  
573 tensive experiments. However, due to space con-  
574 straints, we focus on proposing a knowledge-driven  
575 and dynamic evaluation framework, and do not  
576 delve deeply into methods for optimizing the eval-  
577 uated LLMs based on JudgeAgent. In Section 4  
578 Experiments, we only employ simple prompt injec-  
579 tion optimization as a validation approach to verify  
580 the evaluation’s effectiveness. In future work, we  
581 plan to further explore the potential of this evalu-  
582 ation paradigm, conducting in-depth research on  
583 how to perform targeted training and optimization  
584 of the evaluated LLMs based on JudgeAgent’s eval-  
585 uation results, thereby establishing an efficient and  
586 comprehensive iteration cycle for LLMs.

## 587 Ethical Considerations

588 In this study, all datasets used, including MedQA,  
589 MultiHop-RAG, and QuALITY, were sourced in  
590 compliance with relevant usage guidelines, ensur-  
591 ing no violation of privacy. No personally identi-  
592 fiable information was used, and no experiments  
593 were conducted that could raise privacy or secu-  
594 rity concerns. We are committed to maintaining  
595 transparency and integrity throughout the research  
596 process.

## 597 References

598 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
599 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
600 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
601 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
602 cal report. *arXiv preprint arXiv:2303.08774*.

603 Zachary Ankner, Mansheej Paul, Brandon Cui,  
604 Jonathan D Chang, and Prithviraj Ammanabrolu.  
605 2024. Critique-out-loud reward models. *arXiv*  
606 *preprint arXiv:2408.11791*.

607 Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He,  
608 Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao,  
609 Haozhe Lyu, and 1 others. 2023. Benchmarking founda-  
610 tion models with language-model-as-an-examiner.  
611 *Advances in Neural Information Processing Systems*,  
612 36:78142–78167.

613 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,  
614 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,  
615 Cunxiang Wang, Yidong Wang, and 1 others. 2024.  
616 A survey on evaluation of large language models.  
617 *ACM transactions on intelligent systems and technol-*  
618 *ogy*, 15(3):1–45.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
Ashish Sabharwal, Carissa Schoenick, and Oyvind  
Tafjord. 2018. Think you have solved question an-  
swering? try arc, the ai2 reasoning challenge. *arXiv*  
*preprint arXiv:1803.05457*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, and 1 others. 2021. Training verifiers  
to solve math word problems. *arXiv preprint*  
*arXiv:2110.14168*.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann,  
Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-  
cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and  
1 others. 2025. Gemini 2.5: Pushing the frontier with  
advanced reasoning, multimodality, long context, and  
next generation agentic capabilities. *arXiv preprint*  
*arXiv:2507.06261*.

Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chen-  
hui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu  
Feng, Hanlin Zhao, and 1 others. 2024. Chatglm: A  
family of large language models from glm-130b to  
glm-4 all tools. *arXiv preprint arXiv:2406.12793*.

Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,  
Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan  
Shen, Shengjie Ma, Honghao Liu, and 1 others.  
2024. A survey on llm-as-a-judge. *arXiv preprint*  
*arXiv:2411.15594*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,  
Mantas Mazeika, Dawn Song, and Jacob Steinhardt.  
2020. Measuring massive multitask language under-  
standing. *arXiv preprint arXiv:2009.03300*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,  
Mantas Mazeika, Dawn Song, and Jacob Steinhardt.  
2021. Measuring massive multitask language under-  
standing. In *International Conference on Learning*  
*Representations*.

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei  
Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu,  
Chuancheng Lv, Yikai Zhang, Yao Fu, and 1 oth-  
ers. 2023. C-eval: A multi-level multi-discipline  
chinese evaluation suite for foundation models. *Ad-*  
*vances in Neural Information Processing Systems*,  
36:62991–63010.

Xuhui Jiang, Shengjie Ma, Chengjin Xu, Cehao Yang,  
Liyu Zhang, and Jian Guo. 2025. Synthesize-on-  
graph: Knowledgeable synthetic data generation  
for continue pre-training of large language models.  
*arXiv preprint arXiv:2505.00979*.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,  
Hanyi Fang, and Peter Szolovits. 2021. What disease  
does this patient have? a large-scale open domain  
question answering dataset from medical exams. *Ap-*  
*plied Sciences*, 11(14):6421.

673	Eunsu Kim, Juyoung Suk, Seungone Kim, Niklas Muenighoff, Dongkwan Kim, and Alice Oh. 2025. <b>LLM-as-an-interviewer: Beyond static testing through dynamic LLM evaluation.</b> In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 26456–26493, Vienna, Austria. Association for Computational Linguistics.	<i>Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1733–1747.	729 730 731
674			
675			
676			
677			
678		Yi-Da Tang, Er-Dan Dong, and Wen Gao. 2024. LLMs in medicine: The need for advanced evaluation systems for disruptive technologies. <i>The Innovation</i> , 5(3).	732 733 734
679			
680	Miyoung Ko, Sue Park, Joonsuk Park, and Minjoon Seo. 2024. Hierarchical deconstruction of LLM reasoning: A graph-based framework for analyzing knowledge utilization. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 4995–5027.	Yixuan Tang and Yi Yang. 2024. Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries. <i>arXiv preprint arXiv:2401.15391</i> .	735 736 737
681			
682			
683			
684		Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023a. Mint: Evaluating LLMs in multi-turn interaction with tools and language feedback. In <i>The Twelfth International Conference on Learning Representations</i> .	738 739 740 741 742
685			
686	Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. 2024. Mt-eval: A multi-turn capabilities evaluation benchmark for large language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 20153–20177.	Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, and 1 others. 2023b. Pandalm: An automatic evaluation benchmark for LLM instruction tuning optimization. <i>arXiv preprint arXiv:2306.05087</i> .	743 744 745 746 747 748
687			
688			
689			
690			
691			
692			
693	Yafu Li, Xuyang Hu, Xiaoye Qu, Linjie Li, and Yu Cheng. 2025. Test-time preference optimization: On-the-fly alignment via iterative textual feedback. <i>arXiv preprint arXiv:2501.12895</i> .	Yilei Wang, Jiabao Zhao, Deniz S Ones, Liang He, and Xin Xu. 2025. Evaluating the ability of large language models to emulate personality. <i>Scientific reports</i> , 15(1):519.	749 750 751 752
694			
695			
696			
697	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	753 754 755 756 757
698			
699			
700			
701			
702	Yonatan Oren, Nicole Meister, Niladri S Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. 2023. Proving test set contamination in black-box language models. In <i>The Twelfth International Conference on Learning Representations</i> .	Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yujia Zhou, Wei Shen, Dong Yan, and Yiqun Liu. 2024. Beyond scalar reward model: Learning generative judge from preference data. <i>arXiv preprint arXiv:2410.03742</i> .	758 759 760 761 762
703			
704			
705			
706			
707	Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi, Nikita Nangia, Jason Phang, Angelica Chen, Vishakh Padmakumar, Johnny Ma, Jana Thompson, He He, and 1 others. 2022. Quality: Question answering with long input texts, yes! In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5336–5358.	Zhuohao Yu, Chang Gao, Wenjin Yao, Yidong Wang, Wei Ye, Jindong Wang, Xing Xie, Yue Zhang, and Shikun Zhang. 2024. Kieval: A knowledge-grounded interactive evaluation framework for large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5967–5985.	763 764 765 766 767 768 769
708			
709			
710			
711			
712			
713			
714			
715	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. <i>arXiv preprint arXiv:2311.12022</i> .	J Yuan, P Bao, Z Chen, M Yuan, J Zhao, J Pan, Y Xie, Y Cao, Y Wang, Z Wang, and 1 others. 2023. Advanced prompting as a catalyst: Empowering large language models in the management of gastrointestinal cancers. <i>The Innovation</i> , 5(5).	770 771 772 773 774
716			
717			
718			
719			
720	Rylan Schaeffer. 2023. Pretraining on the test set is all you need. <i>arXiv preprint arXiv:2309.08632</i> .	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging LLM-as-a-judge with mt-bench and chatbot arena. <i>Advances in neural information processing systems</i> , 36:46595–46623.	775 776 777 778 779 780
721			
722	Yi Shi. 2024. Drug development in the ai era: AlphaFold 3 is coming! <i>The Innovation</i> , 5(5).		
723			
724	Zhichao Shi, Shaoling Jing, Yi Cheng, Hao Zhang, Yuanzhuo Wang, Jie Zhang, Huawei Shen, and Xueqi Cheng. 2025. Safetyquizzzer: Timely and dynamic evaluation on the safety of LLMs. In <i>Proceedings of the 2025 Conference of the Nations of the Americas</i>		
725			
726			
727			
728			

## A The Use of Large Language Models

In this paper, we only used LLMs, including DeepSeek-R1 and ChatGPT, for polishing the writing. Specifically, we used LLMs to assist in refining the language, improving readability, and grammar checking. The authors take full responsibility for the content of the paper.

## B Analysis and Proofs of Adaptive Difficulty-Control Mechanisms

In Section 3.1, we have formalized the difficulty control mechanisms as follows:

Let  $[d_1, d_2, d_3]$  respectively denote the difficulty *Easy*, *Medium*, and *Hard*,  $n_{it}, n_{if}$  respectively represent the number of correct and incorrect answers for  $d_i$ -level questions,  $c_{it}$  be the score gain for a correct answer at  $d_i$ , and  $t_{ij}$  be the score threshold between  $d_i$  and  $d_j$ , we can get the average score  $avg_s$  and the difficulty of the next round question generation  $d_{next}$ :

$$avg_s = \frac{\sum_{i=1}^3 c_{it} n_{it}}{\sum_{i=1}^3 (n_{it} + n_{if})}$$

$$d_{next} = \begin{cases} d_1, & avg_s \leq t_{12} \\ d_2, & avg_s \leq t_{23} \\ d_3, & avg_s > t_{23} \end{cases}$$

where  $[c_{1t}, c_{2t}, c_{3t}] = [1, 1.5, 2]$  and  $[t_{12}, t_{23}] = [0.5, 1]$  in our experiments. In this section, we will illustrate the criteria and rules of the score gain  $c_{it}$  and the threshold  $t_{ij}$ , along with the analysis and proofs. For ease of representation, let  $N = \sum_{i=1}^3 (n_{it} + n_{if})$  in the subsequent analysis.

For the selection of score gains and thresholds, we followed simple cognitive principles to establish several rules that a reasonable difficulty control mechanism should satisfy:

**R1. Balance:** Correctly answering a more difficult question should gain a higher score. If the number of correctly answered *Easy* and *Hard* questions is equal, the level should be considered equivalent to the same total number of *Medium* questions. Additionally, the conventional range of the average score should correspond equally to the three difficulty levels in ascending order.

**R2. Generalizability:** The same formula, score gains, and thresholds should apply even to questions without difficulty labels, where difficulty is estimated solely based on the accuracy.

**R3. Improvability:** Correctly answering questions should increase the average score, and there

should be a possibility of exceeding the threshold to advance to the next difficulty level.

**R4. Stability:** If the capability estimation remains at a certain difficulty for a long period, the likelihood of advancing to the next level should gradually decrease.

Following the above rules, we can analyze the mathematical conditions that  $c_{it}$  and  $t_{ij}$  must satisfy. From **R1**, we can get that  $c_{1t} + c_{3t} = 2c_{2t} \wedge c_{1t} < c_{2t} < c_{3t}$ . Since we have set the score for correctly answering a question without a difficulty level as  $c_{2t}$ , the range of  $avg_s$  can be determined as  $[0, c_{2t}]$ , representing the spectrum from answering all questions incorrectly to answering all correctly. Based on the condition about threshold in **R1**, we can get that  $t_{12} = c_{2t}/3$ ,  $t_{23} = 2c_{2t}/3$  in the case that difficulty levels are absent. Due to **R2**, these threshold conditions can be generalized to the case where questions are provided with difficulty levels.

For **R3**, let  $avg_s$  denote the average score after answering  $N$  questions, and  $avg'_s$  denote the average score after answering  $N + 1$  questions. Then **R3** is equivalent to satisfying the following conditions:

$$\mathbf{C1.} \quad avg'_s = \frac{avg_s N + c_{it}}{N + 1} > avg_s,$$

holds true always when  $avg_s \leq t_{(i,i+1)}$

(1)

$$\mathbf{C2.} \quad \exists avg_{s0}, N_0,$$

$$\text{s.t. } avg'_s = \frac{avg_{s0} N_0 + c_{it}}{N_0 + 1} > t_{(i,i+1)}$$

Simplifying Eq 1.C1, we find that  $c_{it} > avg_s$  holds true always when  $avg_s \leq t_{(i,i+1)}$ . Therefore, we can conclude that  $c_{it} > t_{(i,i+1)}$ , meaning  $c_{2t} > c_{1t} > t_{12}$  and  $t_{23} < c_{2t} < c_{3t}$ . Simplifying Eq 1.C2, we can obtain a relationship between  $avg_{s0}$  and  $N_0$ :

$$avg_{s0} > t_{(i,i+1)} - \frac{c_{it} - t_{(i,i+1)}}{N_0} \quad (2)$$

Combining the condition  $c_{it} > t_{(i,i+1)}$  from Eq 1.C1, this relationship indicates that, under the premise of **R3**, advancing to the next difficulty level by answering questions correctly requires a higher  $avg_s$  as  $N$  increases. This description is essentially **R4**: the longer one remains at a specific difficulty level, the harder it becomes to advance to the next level. Therefore, **R3** $\Rightarrow$ **R4** is true.

Finally, **R4** can be stated as follows: given the average score  $a \leq t_{(i,i+1)}$  after  $N$  questions, the

number  $n$  of consecutive questions that must be answered correctly to surpass the level threshold should increase with  $N$ . We can obtain the following inequality:

$$\begin{aligned} \frac{aN + c_{it}n}{N + n} &> t_{(i,i+1)} \\ \Rightarrow (c_{it} - t_{(i,i+1)})n &> (t_{(i,i+1)} - a)N \end{aligned} \quad (3)$$

If **R4** holds, which means  $n$  should increase as  $N$  increases, combined with the inherent condition  $a \leq t_{(i,i+1)}$ , it follows that  $c_{it} > t_{(i,i+1)}$ , from which **R3** can be deduced. Therefore, **R3** $\Rightarrow$ **R4** and **R4** $\Rightarrow$ **R3** hold simultaneously, meaning **R4** and **R3** are equivalent, demonstrating that we have now analyzed all the mathematical conditions that  $c_{it}$  and  $t_{(i,i+1)}$  must satisfy:

$$\begin{aligned} c_{1t} &< c_{2t} < c_{3t} \\ c_{1t} + c_{3t} &= 2c_{2t} \\ t_{12} &= c_{2t}/3, t_{23} = 2c_{2t}/3 \\ c_{it} &> t_{(i,i+1)} \\ \Rightarrow c_{2t} &> c_{1t} > t_{12}, t_{23} < c_{2t} < c_{3t} \end{aligned} \quad (4)$$

Therefore, we first set  $c_{2t} = 1.5$ , and so  $t_{12} = 0.5$ ,  $t_{23} = 1$ . Then, since  $1.5 = c_{2t} > c_{1t} > t_{12} = 0.5$  and  $1 = t_{23} < c_{2t} = 1.5 < c_{3t}$ , we set  $c_{1t} = 1$  for convenience, and so  $c_{3t} = 2c_{2t} - c_{1t} = 2$ .

## C Construction of Context Graph

For the dataset selected during the Benchmark Grading stage, the process of constructing a context graph from its knowledge base (e.g., the set of all reference texts) can be divided into three components: text chunking, node construction, and node linking.

**Text Chunking:** To preserve the integrity of knowledge, the granularity of chunking is not refined to the sentence level. Instead, several sentences or an entire paragraph are grouped to form a single unit, treated as the minimal hierarchical ‘‘chunk’’. On this basis, chunks are further aggregated into higher-level ‘‘documents’’ according to whether they belong to the same article or serve as reference texts for the same question.

**Node Construction:** An LLM is employed to extract entities from each chunk, and the detailed prompt is provided in Appendix G. Each extracted entity serves as the subject for constructing a node. All chunks containing the same entity are assigned to the corresponding entity node. Formally, given

an entity  $e$ , along with the chunks containing the entity  $\{c_1, c_2, \dots\}$  and the documents containing the entity  $\{d_1, d_2, \dots\}$ , the node in the context graph can be defined as follows:

$$\begin{aligned} N &= (e, \mathcal{C}, \mathcal{D}), \\ \text{where } \mathcal{C} &= \{c_1, c_2, \dots\}, \mathcal{D} = \{d_1, d_2, \dots\} \end{aligned} \quad (5)$$

In practice, we store only the IDs of chunks and documents within each node to conserve space.

**Node Linking:** During construction, if two entities appear together in the same chunk or the same document-level texts, their corresponding nodes in the context graph are treated as neighbors. Formally, given two nodes  $N_1 = (e_1, \mathcal{C}_1, \mathcal{D}_1)$  and  $N_2 = (e_2, \mathcal{C}_2, \mathcal{D}_2)$ , if  $\exists c_{1i} \in \mathcal{C}_1, c_{2j} \in \mathcal{C}_2$ , s.t.  $c_{1i} \equiv c_{2j}$  or  $\exists d_{1i} \in \mathcal{D}_1, d_{2j} \in \mathcal{D}_2$ , s.t.  $d_{1i} \equiv d_{2j}$ , then  $N_1$  and  $N_2$  will be treated as neighbors in the context graph.

## D Statistics and Preprocessing of Datasets

### D.1 Details of Datasets

The following benchmarks are used in our experiments, whose details are shown in Table 4.

**MedQA**(Jin et al., 2021) contains multiple-choice questions in the style of the Medical Licensing Examination. Questions in this dataset are collected from medical board exams in the US, Mainland China, and Taiwan, where human doctors are evaluated on their professional knowledge and ability to make clinical decisions. The background knowledge texts of MedQA are provided in the form of additional complete articles, and the questions only provide meta information.

**MultiHop-RAG**(Tang and Yang, 2024) consists of phrase Q&A queries, their ground truth answers, and the associated supported evidence constructed from news articles published between September and December 2023. The background knowledge texts of MultiHop-RAG are provided as supporting evidence along with the questions.

**QuALITY**(Pang et al., 2022) is a multiple-choice question dataset for long document comprehension, whose questions are written and validated by human contributors based on the long passages. The sources of QuALITY include: (1) Project Gutenberg fiction stories, which are mostly science fiction; (2) Slate magazine articles from the Open American National Corpus; (3) other nonfiction articles taken from The Long+Shor, Freesouls, and the book Open Access. QuALITY is organized by articles, with each data item consisting of a long

Datasets	Question Type	Categories	Splits	Used Splits	Split Size
MedQA	multiple-choice	medical clinical	train validation test	test	1273
MultiHopRAG	phrase QA	technology entertainment sports science business health	test	test	2556
QuALITY	multiple-choice	fiction stories magazine articles long articles	train validation test	validation	2086

Table 4: Details of datasets in our experiments. The language of all the datasets is English.

957 article and several related questions, and the arti-  
958 cle serves as the background knowledge for each  
959 question.

## 960 D.2 Preprocessing of Datasets

961 To verify JudgeAgent’s ability to evaluate knowl-  
962 edge deficiencies in target models, we remove  
963 the background knowledge of the questions from  
964 MedQA and MultiHopRAG during evaluation.  
965 Otherwise, it would only test the target’s reading  
966 comprehension ability rather than knowledge defi-  
967 ciencies. In contrast, we provided the background  
968 text when using QuALITY for evaluation, as the  
969 questions in QuALITY are highly dependent on  
970 the text content, and many of the texts are fictional  
971 narratives. Therefore, we use QuALITY to vali-  
972 date the JudgeAgent’s effectiveness in guiding com-  
973 prehension and reasoning rather than discovering  
974 knowledge deficiencies.

975 Additionally, we have reprocessed the difficulty  
976 levels of the questions from QuALITY. Based on  
977 the accuracy of human annotators in answering  
978 questions, QuALITY originally classified ques-  
979 tions into two levels, *Easy* and *hard*, using a 50%  
980 accuracy threshold. To align with the difficulty lev-  
981 els defined by JudgeAgent’s difficulty control mod-  
982 ule (*Easy*, *Medium*, and *Hard*), we re-labeled the  
983 questions using the same accuracy criteria. Ques-  
984 tions of QuALITY with an accuracy below 1/3 are  
985 re-labeled as *Easy*, those with an accuracy between  
986 1/3 and 2/3 as *Medium*, and the rest as *Hard*.

## E Additional Experiment Analysis 987

### E.1 Efficiency of JudgeAgent 988

Target	MedQA		MultiHopRAG	
	avg.token	avg.time	avg.token	avg.time
GLM4-Flash	1558	7.24	2424	23.38
Qwen3	1909	9.95	3453	30.96

Table 5: Efficiency of JudgeAgent. *avg.tokens* and *avg.time* respectively denote the average tokens and time (seconds) cost of JudgeAgent per seed question.

989 We calculated the time and token costs incurred  
990 by JudgeAgent during multi-turn interactions with  
991 GLM4-Flash and Qwen3 as target models, calcu-  
992 lated as averages based on the number of questions  
993 in the benchmark selected during the Benchmark  
994 Grading stage. The results are shown in Table 5.

### E.2 Stability of JudgeAgent 995

996 To investigate the stability of JudgeAgent’s results,  
997 this section presents additional experiments con-  
998 ducted under 5 different random seeds, with the  
999 standard deviation calculated. Experiments are  
1000 conducted using Qwen3 and GLM4-Flash as tar-  
1001 get models on the MedQA and MultiHop-RAG  
1002 benchmarks. The results are shown in Table 6.

1003 The experiment results demonstrate that the re-  
1004 sults across multiple random seeds remain stable  
1005 and are consistent with the ranking of results pre-  
1006 sented in the main results in Table 1. This stability  
1007 benefits from JudgeAgent’s context-graph control  
1008 and knowledge-driven synthesis pathway for eval-  
1009 uation questions. These mechanisms endow the

Target	MedQA			
	ACC1	ACC2	CR↑	CtE↓
Qwen3	91.76±0.09	96.18±0.17	4.83±0.17	0.42±0.08
GLM4-Flash	80.36±0.34	92.65±0.29	13.00±0.38	0.72±0.18

Target	MultiHopRAG			
	ACC1	ACC2	CR↑	CtE↓
Qwen3	63.97±0.53	70.41±0.59	15.78±1.29	9.34±1.35
GLM4-Flash	51.54±0.38	66.56±0.51	23.44±0.88	8.22±1.03

Table 6: The experiment results with standard deviation calculated based on 5 different random seeds.

generated questions and the evaluation feedback with controllability and stability.

### E.3 The Resilience Against Data Contamination

**Can JudgeAgent mitigate the challenges of data contamination in static benchmarking paradigms?** In this section, we simulate a scenario where the static benchmarking evaluation paradigm suffers from data contamination by deliberately exposing the evaluation questions to LLMs during their training process. We selected Llama3-8B-Instruct, Mistral-7B-Instruct-v0.3, and Qwen2.5-7B-Instruct as base models, and constructed supervised fine-tuning (SFT) data from the MedQA and MultiHop-RAG benchmarks, which are intended for evaluation in the main experiments. These training data were used to fine-tune the selected base models. By comparing the performance differences between the original base models and the fine-tuned models on both the static benchmark questions and the extended questions generated by JudgeAgent, we analyze and verify the resilience of JudgeAgent and its derivative JudgeAgent against data contamination.

Specifically, the procedure can be summarized as follows: for a base LLM  $\mathcal{M}$  and a static evaluation benchmark  $\mathcal{D}$ , a fine-tuned LLM  $\mathcal{M}$ -sft is obtained by supervised fine-tuning with the training data constructed from  $\mathcal{D}$ . The performance of  $\mathcal{M}$  and  $\mathcal{M}$ -sft, which is measured by the accuracy (ACC) in answering the questions, is then compared on both  $\mathcal{D}$  and the extended questions  $\mathcal{D}@K$  generated at the  $K$ -th iteration. The severity of data contamination ( $\Delta$ ) is measured by the improvement in performance from the fine-tuned model  $\mathcal{M}$ -sft to the base model  $\mathcal{M}$ , which is formalized as  $\Delta = ACC_{\mathcal{M}\text{-sft}} - ACC_{\mathcal{M}}$ . To mitigate

the effects of the LLM’s randomness, each question was answered by the LLM 5 times. If the LLM produced a correct answer in three or more of these trials, it was considered to have answered the question correctly. We conduct our experiments on an Ubuntu machine with one 40GB NVIDIA A100 GPU. The results are shown in Table 7.

The experiment results indicate that, across various LLMs and benchmarks, fine-tuning with evaluation data leads to a notable improvement in model performance on base questions ( $\Delta$ -base), particularly evident on the MultiHop-RAG benchmark. These findings underscore the risks of data contamination: even when the original model exhibits limited performance on benchmarks, exposure to the benchmark evaluation data can artificially inflate its performance. Consequently, the model’s genuine capabilities may be obscured by overestimated benchmark performance, leading to misapplication in scenarios beyond the actual capabilities.

In contrast, when evaluated on extended questions generated by JudgeAgent, the fine-tuned models and base models show little difference in performance. Additionally, fine-tuning even resulted in a decline on the MedQA benchmark. These results suggest that under the JudgeAgent dynamic evaluation paradigm, the generated questions maintain the validity of the evaluation, even when the original questions have been exposed to the target model. Furthermore, under the same setting of LLM and benchmark, there is a marked gap between the performance gain on base questions and extended questions after fine-tuning, demonstrating that JudgeAgent exhibits considerable resilience to data contamination.

Additionally, we conducted a comparative experiment providing the target models with both standard evaluation feedback from JudgeAgent and

Models	MedQA				MultiHop-RAG			
	base	@1	@2	@3	base	@1	@2	@3
Llama3-8B-Instruct	62.17	69.94	67.72	69.77	48.77	28.30	28.38	29.15
Llama3-sft	82.53	65.59	65.80	65.42	88.45	29.59	29.19	29.79
$\Delta$	<b>20.34</b>	-4.35	-1.92	-4.35	<b>39.69</b>	1.29	0.81	0.65
Mistral-7B-Instruct-v0.3	57.91	59.40	58.76	56.97	59.51	37.75	38.60	37.22
Mistral-sft	62.73	53.18	51.98	51.98	92.17	32.74	32.58	32.62
$\Delta$	<b>4.82</b>	-6.23	-6.78	-4.99	<b>32.66</b>	-5.01	-6.02	-4.60
Qwen2.5-7B-Instruct	85.93	78.04	75.44	77.31	46.31	20.59	21.64	21.48
Qwen2.5-sft	94.20	78.00	75.78	77.27	88.90	26.60	25.56	24.75
$\Delta$	<b>8.27</b>	-0.04	0.34	-0.04	<b>42.59</b>	6.02	3.92	3.27

Table 7: The results for validating the resilience against data contamination. *base* means the performance on base static benchmark questions, and @*K* means extended questions at the *K*-th iteration.

Target	ACC1	ACC <sub>2<sub>d</sub></sub>	ACC <sub>2<sub>f</sub></sub>
Qwen3	91.71	99.18	96.38
GLM4-Flash	80.09	98.37	92.82
GPT-4.1	84.97	98.28	92.44

Table 8: The comparative results on MedQA of directly providing correct answers (ACC<sub>2<sub>d</sub></sub>) and providing feedback from JudgeAgent (ACC<sub>2<sub>f</sub></sub>). ACC1 represents the accuracy of target models before evaluation. All values are percentages.

1085 evaluation feedback containing the correct answers  
1086 as cheating information, investigating the risk of  
1087 data contamination during the evaluation. The  
1088 results are presented in Table 8. As can be ob-  
1089 served, when the correct answers are directly pro-  
1090 vided, nearly all LLMs achieve almost perfect ac-  
1091 curacy. This further illustrates the risk of data con-  
1092 tamination in existing static benchmarks. It also  
1093 demonstrates that the evaluation feedback from  
1094 JudgeAgent does not contain cheating information  
1095 tailored specifically to the seed questions

#### 1096 E.4 Quality of Evaluation Questions

1097 We use Dingo<sup>3</sup> to check the factuality of evaluation  
1098 questions generated by JudgeAgent in multi-turn  
1099 interactions. The results are shown in Table 9,  $N_{eq}$   
1100 represents the number of generated questions in  
1101 Staget Interactive Extension.

1102 As shown in Table 9, in Dingo’s fact-checking  
1103 task, the extended questions from JudgeAgent  
1104 maintain a high fact ratio and detection pass ra-  
1105 tio (using Dingo’s default fact ratio threshold of

<sup>3</sup><https://github.com/MigoXLab/dingo>

Datasets	$N_{eq}$	avg. fact ratio	fact check passed	
			$N_{passed}$	passed ratio
MedQA	3819	97.93%	3680	96.36%
MultiHop-RAG	7668	94.38%	7185	93.70%

Table 9: The fact check results by Dingo.

80%). The average fact ratio is above 94%, while  
the pass ratio exceeds 93% for both MedQA and  
MultiHop-RAG. The results demonstrate that the  
most generated questions exhibit strong factuality,  
ensuring that the dynamic evaluation results are  
based on high-quality questions, and the probabili-  
ty of an evaluated LLM being incorrectly judged  
as wrong when its answer is correct remains low.

As for the diversity of evaluation questions, the  
experiment results of data contamination in Table 7  
in Appendix E.3 can indirectly demonstrate that  
contaminated models perform poorly on synthe-  
sized questions, which indirectly confirms that syn-  
thesized questions differ from seed questions and  
that their diversity is assured.

#### E.5 Cross Validation of the Evaluation Suggestions

**Do the suggestions provided by JudgeAgent only  
take effect for the seed questions?** To address this  
question, we designed cross-validation experiments  
from two perspectives.

First, considering that the evaluation suggestions  
are derived from a comprehensive evaluation of the  
responses to both the seed questions and their ex-  
tended questions, we evaluated and compared the  
accuracy improvement of the suggestions on the  
extended questions versus the seed questions, aim-

ing at verifying that the suggestions do not contain "cheating information" specific to seed questions in the scenario after evaluation. The results are shown in Figure 6.

Secondly, we transferred the evaluation suggestions derived from seed questions to other questions with related knowledge concepts in the benchmark, to verify the effectiveness of the suggestions within the same knowledge domain. Specifically, we categorized questions based on the knowledge entities they contain, formalized as follows: given the context graph  $\mathcal{G}$ , for two questions  $q_1$  and  $q_2$  with knowledge entities  $E_1 = \{e|e \in q_1 \wedge e \in \mathcal{G}\}$  and  $E_2 = \{e|e \in q_2 \wedge e \in \mathcal{G}\}$ , if  $E_1$  and  $E_2$  have a non-empty intersection, then  $q_1$  and  $q_2$  are considered related questions. In this experiment, the suggestions for a question were constructed from the suggestions of its related questions, excluding suggestions from the question itself, to assess and validate the transferability of the suggestions provided by the JudgeAgent. We screened out questions without relevant questions and those with only relevant questions based on non-knowledge entities, such as male, female, 2 years, etc. The results are shown in Table 10.

First, we analyze the difference in the effectiveness of evaluation suggestions on seed questions versus extended questions. As shown in Figure 6, suggestions effectively improve performance for both seed questions and extended questions, with a relatively small gap in the degree of improvement. Notably, when Qwen3 and GPT-4.1 are used as target models, the first round of Qwen3 and the second round of GPT-4.1 exhibit even greater improvement than seed questions. These results indicate that although the suggestions only supplement knowledge for several key concepts, such as the case in Figure 8, such concise suggestions can still benefit both seed and extended questions, demonstrating that JudgeAgent is capable of identifying and addressing knowledge gaps in the target model, rather than simply providing "cheating information" specific to seed questions.

Furthermore, it is observed that the effectiveness of suggestions for third-round questions is consistently low in the experiments. This may be because, by the third round, knowledge path sampling has expanded beyond the scope of knowledge related to seed questions to a broader range. As a result, the generated questions diverge more significantly in core knowledge from earlier questions, thereby reducing the effectiveness of the knowledge guidance

provided in the suggestions.

Next, we analyze the difference in the effectiveness of suggestions on questions that share the same knowledge concepts. As shown in Table 10, compared to their effectiveness on seed questions, the suggestions exhibit a slight decrease when applied to related questions, as indicated by a decline in CR and  $\Delta$ , and an increase in CtE. But the difference is minor, and the improvement remains notable, suggesting that the suggestions can be effectively transferred to other questions involving the same knowledge concepts. This further demonstrates that JudgeAgent can provide suggestions that do not simply serve as "cheating information" specific to seed questions.

However, the slight decline in evaluation effectiveness also indicates that the knowledge guidance provided in the suggestions is not fully aligned with the related questions. This may be because the overlapping entities between these related questions and seed questions do not correspond to core knowledge concepts. For example, suggestions centered on "blood type" may be transferred to a question where "serum" is the core knowledge concept, resulting in a partial mismatch.

The above experiment results demonstrate that the evaluation suggestions provided by JudgeAgent are not only applicable to seed questions but can also be transferred to other questions that share relevant core knowledge concepts.

## E.6 Supplementary Parameter Analysis

**What is the impact of batch size in the *Benchmark Grading* stage on the evaluations of the JudgeAgent?** In *Benchmark Grading* stage, questions are divided into batches to comprehensively assess the target’s capabilities in a base level. These batches are also the basic units for question extension and evaluation feedback. Given a fixed rounds, the batch size is inversely proportional to the number of batches, extended questions, and evaluation suggestions, thereby influencing the time and resource consumption of the entire evaluation process. **Can the batch size be maximized to reduce resource consumption while maintaining the effectiveness of evaluations?** To address this question, we conducted a parameter analysis experiment, examining the evaluation effectiveness and time consumption under different batch sizes. The results are shown in Figure 7.

As observed from the trend of the curves in Figure 7, both the target’s accuracy after receiv-

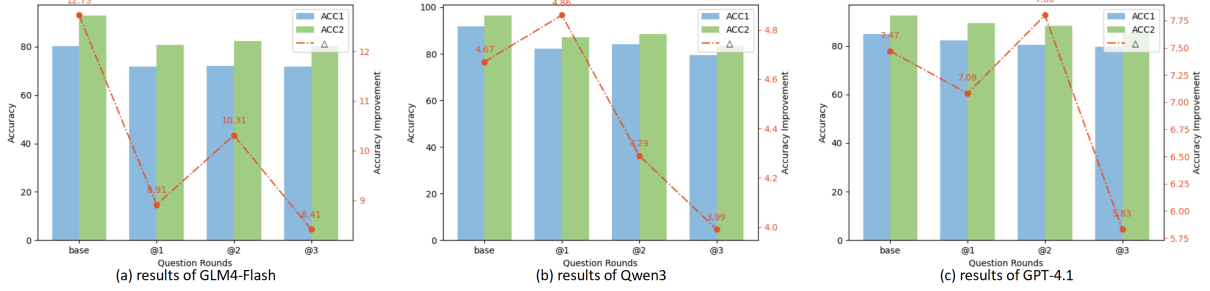


Figure 6: The cross-validation results of extended questions versus seed questions on MedQA. ACC1 and ACC2 indicate the accuracy before and after evaluation.  $\Delta = \text{ACC2} - \text{ACC1}$  refers to the overall accuracy improvement. All the values are percentages. @K represents the results of the questions expanded at the K-th round.

Target	ACC1	Non-transfer				Transfer			
		ACC2	CR $\uparrow$	CtE $\downarrow$	$\Delta \uparrow$	ACC2	CR $\uparrow$	CtE $\downarrow$	$\Delta \uparrow$
GLM4-Flash	78.42	87.49	18.30	9.23	9.07	85.56	16.98	9.84	7.13
GPT-4.1	84.21	91.09	10.23	3.35	6.88	89.78	9.50	3.93	5.57
Qwen3	92.44	96.06	5.96	2.34	3.62	95.17	5.17	2.44	2.73

Table 10: The cross-validation results of transferring suggestions to related questions. All the values are percentages. *Non-transfer* refers to suggestions being applied to the seed questions, whereas *Transfer* refers to them being applied to related questions.

ing evaluation suggestions (ACC2) and the average time consumption per question for evaluation (Time Cost) decrease as the batch size increases. Among them, the decline rate in ACC2 gradually accelerates with larger batch sizes, while the decline rate in Time Cost gradually slows, exhibiting a marginal effect. The reason for the marginal effect in Time Cost lies in the fact that the time required for the target model to answer the seed questions, which is a component of the overall evaluation process, varies little with changes in batch size. As a result, there is a threshold beyond which further reductions in time cost have diminishing returns.

The decline in ACC2 with increasing batch size can be attributed to the expansion of the question’s knowledge domain. As the batch size grows, the generated questions during evaluation become more heterogeneous and less coherent with the knowledge relevant to seed questions, making it difficult for JudgeAgent to identify appropriate knowledge guidance from the dispersed question-answer pairs. Moreover, when the batch size exceeds the number of extension rounds, the disorder of knowledge scopes intensifies more rapidly, ultimately leading to a sharp drop in accuracy. The results in Figure 7 show that when the batch size reaches 9, the evaluation suggestions even become counter-productive ( $\text{ACC2} < \text{ACC1}$ ), interfering with the

normal reasoning of the target model.

Therefore, considering both the evaluation time cost and effectiveness, we selected a batch size of 3 in our experiments as a balanced choice.

## F Algorithm

To clarify the entire workflow of JudgeAgent, we use pseudocode to show the dynamic evaluation process in Algorithm 1, the validation process of evaluation in Algorithm 2, the construction of context graph in Algorithm 3, and the generation of extended questions in Algorithm 4.

## G Prompts

We show the detailed prompt for Entity Extraction in Prompt 11, the prompt for generating extended questions in Prompt 12, the prompt for evaluating the performance of the target LLM in Prompt 13, and the prompt for querying the target LLM with suggestions in Prompt 14.

## H Detailed Content of Case Study

The detailed content of the comparative case in section 4.6 is shown in Figure 8.

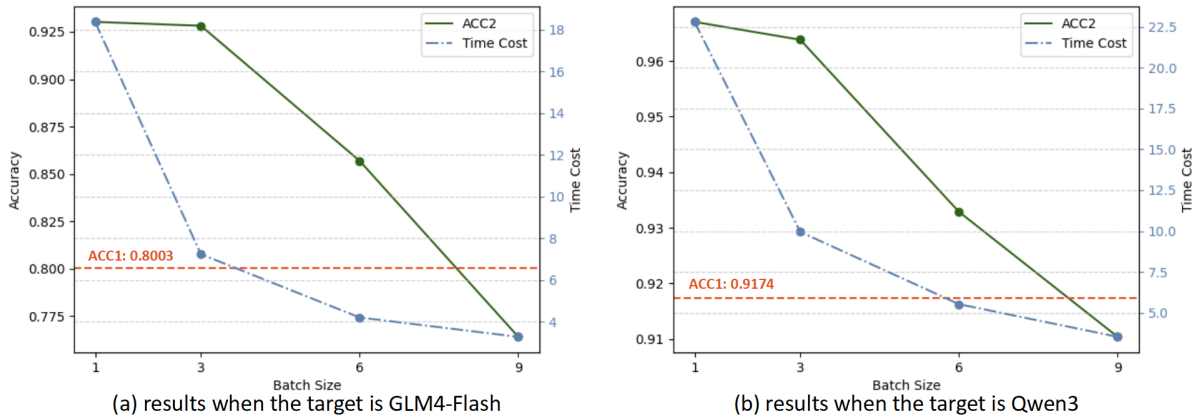


Figure 7: The results of different batch sizes on MedQA. ACC1 and ACC2 indicate the accuracy before and after receiving evaluations. Time Cost is the average time consumption for each question.

---

### Prompt for Entity Extraction

---

Please identify and label the entities in the following multiple sentences, and return the entity labeling results for each sentence.

The results for each sentence should be independent, in JSON format, containing the sentence number, sentence text, and the list of recognized entities (including entity text, type, and position).

Return format is a dictionary, with only one key 'labeled\_data', and the value is a list, each element is a dictionary containing the sentence text and the entity list.

```

{{
  "labeled_data":
  [
    {"text": "Sentence 1", "entity_list": [{"entity_text": "", "entity_type": ""}]},
    {"text": "Sentence 2", "entity_list": [{"entity_text": "", "entity_type": ""}]},
    ...
  ]
}}

```

Notice that "text" should be only the sentence, not the whole article. Sentence list:

```
{ Sentences }
```

---

Table 11: Prompt for entity extraction.

---

### Prompt for Generating Extended Questions

---

As an interviewer, you are tasked with designing questions based on the provided texts. Your role involves crafting questions and correct answers that fulfill the following criteria:

1. **Focus on the Entity**: Ensure all questions consistently center around the specified entity from the article.
2. **Ensure Accuracy and Conciseness of Answers**: Verify that the provided answer is both correct for your designed question within the context and logic of the given text fragments, and ensure the answer is sufficiently concise—presented as a word or phrase, avoiding redundancy.
3. **Conform to difficulty requirements**: You need to design questions for the required difficulty levels, with specific requirements as follows:
  - (1). [easy]: **Encourage Knowledge Memorization**: design questions that assess whether respondents have memorized relevant knowledge. Create questions by directly extracting and blanking out content from the given passage.
  - (2). [medium]: **Encourage Knowledge Comprehension**: Design questions that prompt respondents to dissect and comprehend concepts involved in the topic. Avoid assessing only superficial knowledge retention.
  - (3). [hard]: **Encourage Knowledge Deep Analysis**: Design questions that prompt respondents to engage in deep thinking and analysis. Avoid merely testing knowledge recall or conceptual comprehension; do not simply extract fragments from the given passage to create fill-in-the-blank items. Encourage respondents to focus on entities within the question and employ logical skills for complex reasoning.

Here are examples:

{ *Examples* }

Now, given the following text fragments:

{ *context* }

Based on the provided texts, please generate questions by following the requirements above and referencing the examples.

Output in the specified JSON format below: {{

```
"generated_question":
[
  {{
    "question": "Generated Question",
    "answer": "Correct Answer of Generated Question"
  }},
  ...
]
```

}}

---

Table 12: Prompt for generating extended questions based on different difficulties.

---

**Prompt for Evaluation**

---

The following is the performance of an LLM in answering a series of questions:

{ *list of questions, correct answers, and LLM's answers* }

Please evaluate and analyze the interviewee's performance based on the above performance using concise language from the following perspectives, and provide suggestions that help the LLM answer the same questions better. Suggestions should provide specific and detailed guidance on logical thinking steps, required knowledge, and abilities, ensuring the LLM can answer correctly for the same questions.

Output in the following JSON format:

```
{{  
  "flaws_knowledge": "The lack of background knowledge.",  
  "flaws_capability": "The flaws in logic and capability.",  
  "comprehensive_performance": "The Comprehensive performance of all questions.",  
  "suggestions": "Suggestions that help the LLM answer questions better"  
}}
```

---

Table 13: Prompt for evaluating the performance of the target LLM.

---

**Prompt for Querying LLM with suggestions**

---

Please complete the following question:

[question]: { *question* }

In your previous responses to these questions, the interviewer has provided the following suggestions for you to help you answer better:

[suggestions]: { *suggestions* }

Please consider the above [suggestions], and answer the above [question], in the following JSON format:

```
{{ "answer": "Your answer" }}
```

---

Table 14: Prompt for querying the target LLM with suggestions.

### Base Question

In *closed abdominal injury*, rupture of which of the following sites results in the latest onset of *peritonitis* symptoms?  
A. Colon.      B. Duodenal bulb.      C. Ileum.      D. Jejunum

A

### Direct Evaluation

LLM Evaluator

To improve, LLM should understand following knowledge:  
(1) **Colon rupture**, due to the thick consistency and slower leakage rate, and with parts of the colon being retroperitoneal, **peritonitis symptoms may be delayed**.  
(2) **Duodenal bulb rupture** causes leakage of pancreatic juice and bile into the retroperitoneal space rather than directly into the peritoneal cavity. Thus, **early symptoms are subtle, and peritonitis signs (such as tenderness and guarding) appear later**.  
(3) **Ileum rupture** results in **rapid leakage** of intestinal contents into the peritoneal cavity, causing **acute peritonitis**.  
(4) **Jejunum rupture** is similar to ileum rupture, with contents entering the peritoneal cavity directly, leading to **rapidly developing peritonitis**.  
**Suggestions:** Familiarity with **typical peritonitis symptoms** (e.g., pain, tenderness, guarding) and **their timing**, as well as differences in clinical presentations of common abdominal injuries.

In *closed abdominal injury*, rupture of which of the following sites results in the latest onset of *peritonitis* symptoms?  
A. Colon.      B. Duodenal bulb.      C. Ileum.      D. Jejunum

A

### Agent-as-Interviewer

Entities: [ *closed abdominal injury*, *peritonitis* ]

In *closed abdominal injuries*, **splenic rupture** accounts for 20%-40% of cases, while in open abdominal injuries, it constitutes about 10%. A spleen with chronic pathologies (such as schistosomiasis, malaria, lymphoma, etc.) is more prone to rupture. Based on pathological anatomy, splenic rupture can be classified into three types: central rupture (located deep within the splenic parenchyma), subcapsular rupture (located at the periphery of the splenic parenchyma), and true rupture (involving the capsule).

**Duodenal injury**, if it occurs in the intraperitoneal portion, can cause pancreatic juice and bile to leak into the abdominal cavity, leading to early symptoms of **peritonitis**. Although preoperative diagnosis may not clearly pinpoint the site of injury, the obvious symptoms generally do not delay the timing of surgery. In contrast, duodenal rupture caused by **closed abdominal injury** often presents with subtle early signs, making timely recognition difficult. Vigilance should be heightened in the following situations: persistent and progressively worsening pain in the right upper abdomen or flank, which may radiate to the right shoulder and right testis; obvious localized tenderness in the right upper abdomen and flank; relatively mild abdominal signs accompanied by deteriorating general condition; occasional vomiting of blood-tinged material; elevated serum amylase; abdominal X-ray showing blurred outlines of the psoas muscle, sometimes with mottled changes (pneumoretroperitoneum) that gradually spread; extravasation observed when water-soluble contrast is introduced via a gastric tube; CT or MRI revealing air bubbles in the retroperitoneum and anterior to the right kidney; and on digital rectal examination, the presence of crepitus in the anterior rectal wall, indicating that gas has reached the pelvic retroperitoneal space.

**Q1:** A patient with *closed abdominal injury* presents with right upper quadrant pain, bloody vomitus, elevated serum amylase, and retroperitoneal air accumulation. What is the most likely cause?  
A. Splenic rupture.    B. Hepatic rupture.  
C. Gastric rupture.    D. **Duodenal rupture.**

C

**Q2:** A spleen with chronic pathological changes is more prone to rupture. Which of the following diseases most commonly increases the risk of *splenic rupture*?  
A. Diabetes.      B. **Schistosomiasis.**  
C. Hypertension.    D. Asthma.

B

Based on the overall performance, the LLM **lacks understanding of the symptoms associated with duodenal rupture caused by closed abdominal injury**. Thus, the LLM should understand the knowledge: peritonitis develops more rapidly when organs containing bacteria (such as colon, ileum) rupture, exposing the sterile peritoneal cavity to contamination. Duodenal injury, if it occurs in the intraperitoneal portion, can cause pancreatic juice and bile to leak into the abdominal cavity, leading to early symptoms of peritonitis. Although preoperative diagnosis may not clearly pinpoint the site of injury, the obvious symptoms generally do not delay the timing of surgery. **In contrast, duodenal rupture caused by closed abdominal injury often presents with subtle early signs, making timely recognition difficult.** The most delayed peritonitis symptoms after closed abdominal trauma are with duodenal bulb injury, due to the less severe initial bacterial contamination. This knowledge should be referenced when addressing similar clinical scenarios about the timing of peritonitis onset in relation to anatomical sites.

In *closed abdominal injury*, rupture of which of the following sites results in the latest onset of *peritonitis* symptoms?  
A. Colon.      B. **Duodenal bulb.**      C. Ileum.      D. Jejunum

B

Figure 8: The detailed content of the case.

---

**Algorithm 1** Dynamic evaluation process of JudgeAgent

---

```
1: Input: Target LLM  $\mathcal{M}_t$ , Base dataset  $\mathcal{D} = \{(q_i, a_i, d_i)\}_{i=1}^N$  (each item include question  $q$ , answer  $a$ , and difficulty  $d$ ), Knowledge bases  $\mathcal{K} = \{k_1, k_2, \dots\}$ , Core LLM of JudgeAgent  $\mathcal{M}_c$ , Predefined batch capacity  $N_B$ , Max extension round  $RND_e$ , Max hop of sampling  $H$ 
2: Output: Evaluation Score  $sc$ , Batched dataset with suggestions  $\mathcal{D}_S = \{(\mathcal{B}_1, s_1), \dots\}$ 
   // Construct context graph
3:  $\mathcal{G} \leftarrow \text{CONSTRUCT\_CONTEXT\_GRAPH}(\mathcal{K}, \mathcal{M}_c)$ 
   // Split base dataset into batches
4:  $\{\mathcal{B}_1 = \{(q_{1i}, a_{1i}, d_{1i})\}_{i=1}^{N_B}, \dots\} \leftarrow \text{SPLIT\_BATCHES}(\mathcal{D}, N_B)$ 
   // Begin Evaluation
5:  $sc \leftarrow 0, N_{total} \leftarrow 0$ 
6:  $\mathcal{D}_S \leftarrow \{\}$ 
7: for  $\mathcal{B} \leftarrow \{\mathcal{B}_1, \mathcal{B}_2, \dots\}$  do
8:    $sc_{\mathcal{B}} \leftarrow 0$ 
9:    $RND_{total} \leftarrow 0$ 
10:   $\mathcal{Q}_{tested} \leftarrow \{\}$ 
   // Stage1: Benchmark Grading
11:  for  $(q, a, d) \leftarrow \mathcal{B}$  do
12:    Get answer  $a_{\mathcal{M}} \leftarrow \text{QUERY\_LLM}(q, \mathcal{M}_t)$ 
13:    if  $a_{\mathcal{M}}$  is correct based on  $q$  and  $a$  then
14:       $sc_{\mathcal{B}} \leftarrow sc_{\mathcal{B}} + \text{DIFFICULTY\_SCORE}(d)$ 
15:    end if
16:     $RND_{total} \leftarrow RND_{total} + 1$ 
17:    Add  $(q, a, d, a_{\mathcal{M}})$  to  $\mathcal{Q}_{tested}$ 
18:  end for
19:  Decide difficulty  $d_e \leftarrow \text{DECIDE\_DIFFICULTY}(sc_{\mathcal{B}}, RND_{total})$ 
   // Stage2: Interactive Extension
20:  for  $i \leftarrow \{1, 2, \dots, RND_e\}$  do
21:     $(q_e, a_e, t) \leftarrow \text{GENERATE\_EXTENDED\_QUESTIONS}(\mathcal{B}, \mathcal{G}, \mathcal{M}_c, H, d_e)$ 
22:    Get answer  $a_{\mathcal{M}} \leftarrow \text{QUERY\_LLM}(q_e, \mathcal{M}_t)$ 
23:    if  $a_{\mathcal{M}}$  is correct based on  $q_e$  and  $a_e$  then
24:       $sc_{\mathcal{B}} \leftarrow sc_{\mathcal{B}} + \text{DIFFICULTY\_SCORE}(d_e)$ 
25:    end if
26:     $RND_{total} \leftarrow RND_{total} + 1$ 
27:    Add  $(q_e, a_e, d_e, a_{\mathcal{M}})$  to  $\mathcal{Q}_{tested}$ 
28:    Decide difficulty  $d_e \leftarrow \text{DECIDE\_DIFFICULTY}(sc_{\mathcal{B}}, RND_{total})$ 
29:  end for
   // Stage3: Evaluation Feedback
30:  Get evaluation suggestions  $s \leftarrow \text{EVALUATE}(\mathcal{Q}_{tested}, \mathcal{M}_c)$ 
31:  Add  $(\mathcal{B}, s)$  to  $\mathcal{D}_S$ 
32:   $sc \leftarrow sc + sc_{\mathcal{B}}$ 
33:   $N_{total} \leftarrow N_{total} + RND_{total}$ 
34: end for
35:  $sc \leftarrow sc / N_{total}$ 
36: return  $sc, \mathcal{D}_S$ 
```

---

---

**Algorithm 2** Validation process of evaluation results from JudgeAgent

---

```
1: Input: Target LLM  $\mathcal{M}_t$ , Batched dataset with suggestions  $\mathcal{D}_S = \{(\mathcal{B}_1, s_1), \dots\}$ , in which  $\mathcal{B}_i = \{(q_{i1}, a_{i1}, d_{i1}), \dots\}$  is a batch of base dataset  $\mathcal{D}$ , and  $\mathcal{S}_i = \{s_{i1}, \dots\}$  is relevant suggestions from JudgeAgent.
2: Output: Accuracy of target LLM before evaluation  $acc_1$ , Accuracy after evaluations  $acc_2$ , Correction Rate  $cr$ , Correct-to-Error Rate  $ce$ 
   // Initialize counter of questions
3:  $N_{acc1} \leftarrow 0$ ,  $N_{acc2} \leftarrow 0$ ,  $N_{cr} \leftarrow 0$ ,  $N_{ce} \leftarrow 0$ ,  $N_{total} \leftarrow 0$ 
4: for  $(\mathcal{B}, s) \leftarrow \mathcal{D}_S$  do
5:    $N_{\mathcal{B}} \leftarrow \text{LEN}(\mathcal{B})$ 
6:    $N_{total} \leftarrow N_{total} + N_{\mathcal{B}}$ 
7:   for  $i \leftarrow \{1, 2, \dots, N_{\mathcal{B}}\}$  do
8:      $(q, a, d) \leftarrow \mathcal{B}[i]$ 
9:     Get answer1  $a_1 \leftarrow \text{QUERY\_LLM}(q, \mathcal{M}_t)$ 
10:    Get answer2  $a_2 \leftarrow \text{QUERY\_LLM\_WITH\_SUGGESTIONS}(q, s, \mathcal{M}_t)$ 
11:    correct1  $\leftarrow$  whether  $a_1$  is correct based on  $q$  and  $a$ 
12:    correct2  $\leftarrow$  whether  $a_2$  is correct based on  $q$  and  $a$ 
13:    if correct1 then
14:       $N_{acc1} \leftarrow N_{acc1} + 1$ 
15:      if not correct2 then
16:         $N_{ce} \leftarrow N_{ce} + 1$ 
17:      end if
18:    end if
19:    if correct2 then State  $N_{acc2} \leftarrow N_{acc2} + 1$ 
20:      if not correct1 then
21:         $N_{cr} \leftarrow N_{cr} + 1$ 
22:      end if
23:    end if
24:  end for
25: end for
26:  $acc_1 \leftarrow N_{acc1}/N_{total}$ ,  $acc_2 \leftarrow N_{acc2}/N_{total}$ ,  $cr \leftarrow N_{cr}/N_{total}$ ,  $ce \leftarrow N_{ce}/N_{total}$ 
27: return  $acc_1$ ,  $acc_2$ ,  $cr$ ,  $ce$ 
```

---

---

**Algorithm 3** Construction process of context graph

---

```
1: Input: Knowledge Base of Dataset  $\mathcal{K} = \{k_1, k_2, \dots\}$ , an LLM  $\mathcal{M}$ 
2: Output: Context Graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ , in which  $\mathcal{N}$  is the node set, and  $\mathcal{E}$  is the edge set.
3:  $\mathcal{N} \leftarrow \{\}$ ,  $\mathcal{E} \leftarrow \{\}$  // Initialize node set and edge set as empty dictionary
4: for  $k \leftarrow \mathcal{K}$  do
    // Splitting text into chunks
5:    $C = \{c_1, c_2, \dots\} \leftarrow \text{SPLIT\_TO\_CHUNKS}(k)$ 
    // Prompting LLM to label entities
6:    $P_e \leftarrow \emptyset$ 
7:   for  $c \leftarrow \{c_1, c_2, \dots\}$  do
8:      $S_e = \{e_1, e_2, \dots\} \leftarrow \text{LABEL\_ENTITY}(c, \mathcal{M})$ 
9:     for  $e \leftarrow S_e$  do
10:      if  $e$  not in  $\mathcal{N}$  then
11:         $\mathcal{N}[e] \leftarrow (e, \mathcal{C} = \emptyset, \mathcal{D} = \emptyset)$ 
12:         $\mathcal{E}[e] \leftarrow \emptyset$ 
13:      end if
14:      Add  $c$  to set  $\mathcal{N}[e].\mathcal{C}$ 
15:      Add  $k$  to set  $\mathcal{N}[e].\mathcal{D}$ 
16:      Expand set  $\mathcal{E}[e]$  with  $S_e$ 
17:    end for
18:    Expand set  $P_e$  with  $S_e$ 
19:  end for
20: for  $e \leftarrow P_e$  do
21:   Expand set  $\mathcal{E}[e]$  with  $P_e$ 
22: end for
23: end for
24:  $\mathcal{G} \leftarrow (\mathcal{N}, \mathcal{E})$ 
25: return  $\mathcal{G}$ 
```

---

---

**Algorithm 4** Generation process of extended questions

---

1: **Input:** Base Question  $\mathcal{Q} = \{q_1, q_2, \dots\}$ , Context Graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ , LLM  $\mathcal{M}$ , Max hop of path  $H$ , Difficulty  $d$

2: **Output:** Extended question with its answer and background text  $(q_e, a_e, t)$   
*// Prompting LLM to label entities from  $\mathcal{Q}$*

3:  $S_e \leftarrow \{\}$

4: **for**  $q \leftarrow \mathcal{Q}$  **do**

5:      $set_e = \{e_1, e_2, \dots\} \leftarrow \text{LABEL\_ENTITY}(q, \mathcal{M})$

6:     Add  $set_e$  to  $S_e$

7: **end for**

8:  $e \leftarrow \text{RANDOM\_SAMPLE}(\{e_1, e_2, \dots\}, 1)$   
*// Sample knowledge paths*

9:  $e' \leftarrow$  the most similar entity in  $\mathcal{N}$  of  $\mathcal{G}$

10:  $E_{visited} \leftarrow \{e'\}$

11:  $t \leftarrow \text{RANDOM\_SAMPLE}(\mathcal{N}[e'].C, 1)$

12: **for**  $i \leftarrow \{1, 2, \dots, H - 1\}$  **do**

13:      $E_{candidate} \leftarrow$  the most similar 5 entities to  $e'$  in  $\mathcal{E}[e']$  that not in  $E_{visited}$

14:      $e' \leftarrow \text{RANDOM\_SAMPLE}(E_{candidate}, 1)$

15:      $C_{candidate} \leftarrow$  the most similar 5 chunks to  $q$  in  $\mathcal{N}[e'].C$

16:      $c \leftarrow \text{RANDOM\_SAMPLE}(C_{candidate}, 1)$

17:     Concatenate  $c$  to new line of  $t$

18:     Add  $e'$  to  $E_{visited}$

19: **end for**

20:  $(q_e, a_e) \leftarrow \text{GENERATE\_QUESTION}(t, \mathcal{M})$

21: **return**  $(q_e, a_e, t)$

---