

# FACTGUARD: DETECTING UNANSWERABLE QUESTIONS IN LONG-CONTEXT TEXTS FOR RELIABLE LLM RESPONSES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models (LLMs) have demonstrated significant advances in reading comprehension. However, a persistent challenge lies in ensuring these models maintain high accuracy in answering questions while reliably recognizing unanswerable queries. This issue remains critical, particularly as the length of supported contexts continues to expand. To address this challenge, we propose a collaborative multi-task workflow called FactGuard to automatically generate evidence-based question-answer pairs and systematically construct unanswerable questions. Using this methodology, we developed the FactGuard-Bench dataset, which comprises 25,220 examples of both answerable and unanswerable question scenarios, with context lengths ranging from 4K to 128K. Experimental evaluations conducted on nine popular LLMs reveal that all LLMs exhibit significant performance gap between answerable and unanswerable questions and the most advanced models achieve only 67.67% overall accuracy. After training with FactGuard-Bench, the model achieves an overall accuracy of 81.17%, along with enhanced reasoning capabilities on unanswerable questions. Our code is publicly available at <https://anonymous.4open.science/r/FACTGUARD-5BBC>

## 1 INTRODUCTION

Comprehending text and answering questions are foundational capabilities in the field of Natural Language Processing (NLP). Over the years, large language models (LLMs) have made substantial progress in reading comprehension, including the ability to process long-context inputs of up to 128K tokens (Yang et al., 2025; Liu et al., 2024). However, LLMs often tend to be overconfident (Slobodkin et al., 2023) and specially face an increased risk of generating hallucination or plausible content on unanswerable questions (Deng et al., 2024). This will undermine confidence in LLM capabilities and diminish their overall reliability.

Extracting answers to answerable questions or providing justifications for why certain questions are unanswerable is equally essential for enhancing the practicality of LLMs. Answerable questions can be resolved using information contained in the provided context, while unanswerable questions arise when the context lacks sufficient or reliable evidence to support a definitive response. Handling unanswerable questions presents a particularly challenging scenario, as it requires LLMs to deeply comprehend the context, accurately determine that the question cannot be answered, and provide appropriate reasons to convince the user.

Recently, many advanced works have made a lot of efforts on unanswerable questions (Deng et al., 2024; Yehuda et al., 2024; Rajpurkar et al., 2018). SQuAD 2.0 (Rajpurkar et al., 2018) focuses on the reading comprehension of models in both answerable and unanswerable questions with manual annotation. Its texts are constrained by a context length of under 4K tokens and it does not include explicit refusal responses for unanswerable questions. SelfAware (Yin et al., 2023) employs a straightforward approach that prompts LLMs to detect unanswerable questions and response to them using predefined replies such as, "The answer is unknown". KUQ (Amayuelas et al., 2024) handles open-source LLMs on Known-Unknown questions in open-ended question-answering scenarios rather than questions related to reading comprehension. Self-Aligned method (Deng et al., 2024)

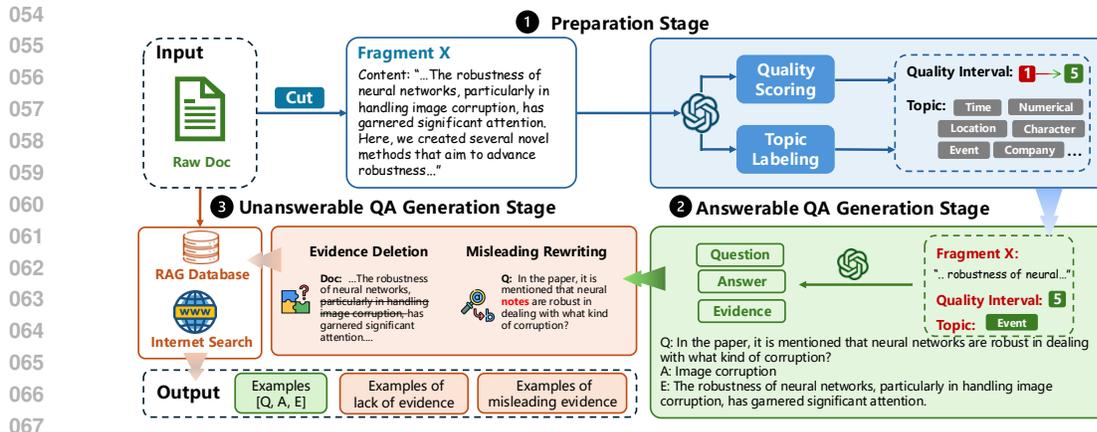


Figure 1: Illustration of FactGuard for data synthesis in a collaborative multi-task workflow framework.

mainly focuses on reasoning responses to unanswerable questions and does not pay attention to long context and it also requires manually labeled questions as seed data.

To overcome the above limitations, we propose a novel approach that employs a collaborative multi-task workflow framework to enable automated data augmentation. We introduce **FactGuard-Bench**, a reading comprehension dataset of 25,220 questions (8,829 answerable and 16,391 unanswerable) with context lengths ranging from 4K to 128K, developed through our framework. Experiments show that even the best-performing model achieves an overall accuracy of 67.67%, with significantly lower performance on unanswerable questions compared to answerable ones. Through further training, we achieved an overall accuracy of 81.17%, along with enhanced reasoning capabilities on unanswerable questions.

We highlight our contributions as follows:

1. **Innovative collaborative multi-task workflow for Data Augmentation:** We introduce **FactGuard**, a collaborative multi-task workflow framework designed to dynamically generate both answerable and unanswerable questions through a coordinated multi-step process. This approach produces contextually challenging examples that enhance the comprehension ability of LLMs.
2. **Development of Benchmark for Long-Context Evaluation:** We curate **FactGuard-Bench**, a long texts benchmark specifically tailored to assess the ability of LLMs to handle answerable and unanswerable questions.
3. **Enhancement of LLMs on unanswerable questions:** Experiments with state-of-the-art LLMs show our method enhanced the model’s ability to handle unanswerable questions and generated well-reasoned answers when solving unanswerable questions.

## 2 RELATED WORK

### 2.1 MACHINE READING COMPREHENSION

Machine reading comprehension (MRC) is a hot research topic in the field of NLP, which focuses on reading documents and answering related questions (Liu et al., 2019; Baradaran et al., 2022). Over the years, machine reading comprehension has garnered significant attention from both academia and industry (Hermann et al., 2015; Liu et al., 2019). With the rapid advancements of large language models (LLMs) Zhao et al. (2023); Liu et al. (2023), retrieval-augmented generation (RAG) has emerged as a promising framework for tackling reading comprehension tasks across diverse specialized domains (Zhao et al., 2024; Lewis et al., 2020). Nevertheless, even state-of-the-art RAG frameworks are susceptible to retrieval accuracy limitations (Hu et al., 2019; Wang et al., 2025), which emphasizes the critical importance of facticity Jacovi et al. (2025); Bi et al. (2025), i.e., the

108 ability of a model to generate factually consistent and verifiable responses in information-seeking  
109 scenarios. In this work, we emphasize scalable and robust evaluation of answerable and unanswerable  
110 questions in reading comprehension.

## 112 2.2 LONG CONTEXT LLMs AND BENCHMARKS

113  
114 Recent studies have emphasized the importance of extending positional embeddings to improve the  
115 ability of LLMs to handle long contexts effectively (Su et al., 2024; Press et al., 2021; Chi et al., 2022).  
116 Closed-source LLMs, in particular, have emerged as leaders in long-context modeling, benefiting  
117 from progressively larger context windows. For instance, models such as GPT-4 (Achiam et al.,  
118 2023) and Gemini Pro 1.5-1000k (Team et al., 2024) are capable of processing increasingly longer  
119 documents, with context lengths ranging from 128k to 1000k tokens. Similarly, open-source LLMs,  
120 including Qwen 2.5 (Yang et al., 2024a) and DeepSeek (DeepSeek-AI, 2024), also support context  
121 lengths of at least 128k tokens. Key benchmarks for assessing long-context capabilities include NIAH  
122 (gkamradt, 2023; Yu et al., 2025), Longbench Series (Bai et al., 2024; 2025), LooGLE (Li et al.,  
123 2024), and L-Eval (An et al., 2024), among others. In FactGuard-Bench, we utilize a wider range of  
124 context lengths to evaluate the LLM’s ability to understand, learn, and reason about information in  
125 text.

## 126 2.3 DETECTION OF UNANSWERABLE QUESTIONS

127  
128 In recent years, studies have increasingly focused on enhancing the ability of reading comprehension  
129 models to detect unanswerable questions. Approaches such as those by (Yin et al., 2023) and  
130 (Slobodkin et al., 2023) employ prompt engineering—for instance, by incorporating hints such as,  
131 “If the question cannot be answered based on the passage, reply ‘unanswerable’”—to improve the  
132 model’s ability to detect unanswerable questions. On the other hand, some methods (Agarwal et al.,  
133 2023; Deng et al., 2024; Rajpurkar et al., 2018) construct datasets related to unanswer questions to  
134 evaluate the model’s ability of detection of unanswerable questions. For example, (Agarwal et al.,  
135 2023) categorized unanswerable questions into five distinct types: Incomplete Information, Future  
136 Questions, Incorrect Information, Ambiguous, and Unmeasurable. They also introduced QnotA—a  
137 dataset consisting of 400 samples designed to support this taxonomy. However, these datasets are  
138 often small in scale, require expensive annotation manpower, and have a short context information.  
139 We automatically constructed FactGuard-Bench generating by LLMs, a large-scale dataset comprising  
140 tens of thousands of long texts. This dataset enables a comprehensive multi-dimensional evaluation  
141 of model capabilities in detecting unanswerable questions.

## 142 3 FACTGUARD METHODOLOGY

143  
144 As shown in Figure 1, we propose FactGuard, a collaborative multi-task workflow framework for  
145 automated data synthesis. FactGuard consists of three primary stages: Preparation Stage, Answerable  
146 QA Generation Stage, and Unanswerable QA Generation Stagee.

### 148 3.1 PREPARATION STAGE

149  
150 We slice the original long document into multiple short text fragments. The window size is kept at  
151 [500, 1000] and slicing is done on a paragraph by paragraph basis. We randomly select Fragment X  
152 for the following sub-steps:

- 153 • **Quality Scoring:** Using LLM, we evaluate Fragment X in terms of fluency, coherence, and  
154 logicity, assigning a quality score on a 5-point scale ( $score_i \in [1, 5]$ ). Fragments with  
155 score lower than 4 will be discarded to ensure their high quality.
- 156 • **Topic Labeling:** Then, we utilize LLM to extract structured information as topic labeling  
157 (e.g., temporal expressions, numerical values, entity, locations, organizations, and events)  
158 from Fragment X. Fragments without clear structured information will be discarded because  
159 these structured information are important for QA generation.

160  
161 After preparation stage, We obtain many high-quality fragments with clear structural information  
from the original long document.

Reasoning	Description	Example
Lack of Evidence	The question is related to the article, but the factual basis is deleted.	<b>Fragment:</b> ...There had been a lack of confidence in Murray since Romani, and the two failed Gaza battles increased his unpopularity among both the infantry and the mounted troops. <del>After the war Allenby acknowledged Murray's achievements in a June 1919 despatch in which he summed up his campaigns...</del> <b>Question:</b> According to this article, in what year did Allenby recognize Murray's accomplishments in his circular? <b>Answer:</b> The question cannot be answered. The article mentions Murray's performance in the battle, but does not mention what year Allenby recognized his accomplishments.
Misleading Evidence	The key information of the question is misaligned against the facts of the article.	<b>Fragment:</b> <a href="#">Global and Local Mixture Consistency Cumulative Learning (GLMC)</a> for Long-Tailed Visual Recognition...The paper introduces GLMC, a one-stage training strategy designed to improve long-tailed visual recognition by enhancing the robustness of the feature extractor and reducing the bias of the classifier towards head classes. GLMC uses a global and local mixture consistency loss and a cumulative head-tail soft label reweighted loss... <b>Raw Question:</b> What are the core ideas behind the <a href="#">Global and Local Mixture Consistency cumulative learning (GLMC)</a> framework and how does it improve long-tailed visual recognition? <b>Question1 with entity substitutions:</b> What are the core ideas behind the <a href="#">Global and Local Augmentation Consistency Learning (GLACL)</a> framework and how does it improve long-tailed visual recognition? <b>Answer1:</b> The article focuses on GLMC and does not mention GLACL. The core ideas of GLACL cannot be answered, but about GLMC... <b>Question2 with impossible condition insertions:</b> What are the core ideas behind the <a href="#">Global and Local Mixture Consistency cumulative learning (GLMC)</a> and framework and how does it improve long-tailed visual recognition <a href="#">on CIFAR-100-LT</a> ? <b>Answer2:</b> The article does not mention CIFAR-100-LT. The question of how GLMC improves long-tailed visual recognition on CIFAR-100-LT cannot be answered, but the article mentioned GLMC improve long-tailed visual recognition by enhancing ...

Table 1: A detailed categorization of unanswerable examples in FactGuard-Bench.

### 3.2 ANSWERABLE QA GENERATION STAGE

On answerable QA generation stage, we generate questions, answers and evidence based on high-quality fragments obtained in preparation stage. Note that evidence consists of specific text segments from fragments that substantiate the answer. This design ensures that each question is firmly grounded in the original long document. Since there are low-quality results for LLM generation, such as questions that are not fluent or the evidence does not come from the fragments, we filter them with quality judgment after answerable QA generation.

After QA Generation stage, we can obtain the answerable questions, answers and evidence derived from the original text.

### 3.3 UNANSWERABLE QA GENERATION STAGE

On unanswerable QA generation stage, we generate unanswerable questions and their corresponding answers based on the answerable questions that have already been generated in the QA generation stage. There are mainly two methods for generating unanswerable QA:

- **Unanswerable questions of lacking evidence:** We simply remove the evidence from fragment, thus making the question unanswerable due to lack of information. For the rejection response, we ask the LLM to provide a reasonable rejection response that echoes the question, and then introduce the main content of the document to prove that the answer cannot be found in the text.
- **Unanswerable questions of misleading evidence:** To create misleading questions, we use LLM to rewrite question through entity substitutions and impossible condition insertions. When rewriting the question through entity replacement, we require that in the rejection responses generated by LLM, it should be indicated that the content appearing in the article is related to the entity before replacement, rather than that of the entity after replacement. When rewriting the question through impossible condition insertions, We require LLM

216 to first refer to the explanation in the rejection response to clarify that the answers to the  
217 questions with impossible condition insertions cannot be found in the original text, and then  
218 answer the original questions before rewriting.  
219

220 As shown in Table 1, a detailed overview of unanswerable examples in FactGuard-Bench can be  
221 found. For unanswerable questions of lacking evidence, we remove the evidence from the original  
222 fragment. For unanswerable questions with misleading evidence, the Fragment remains unchanged,  
223 but we rewrite the questions using entity substitution or impossible condition insertions.  
224

225 After unanswerable QA generation stage, we can obtain the unanswerable questions along with  
226 reasonable response that remain highly relevant to the original text.

227 To ensure the quality of answerable and unanswerable questions, we review process for the generated  
228 data by employing Retrieval Augmented Generation (RAG) techniques. This approach allows  
229 us to extract the top N relevant passages from a lengthy article for short-reading comprehension  
230 and to filter out data that contain conflicting answers. Furthermore, we employ the World Wide  
231 Web to filter common-sense knowledge, effectively circumventing the inherent conflict between  
232 context-faithfulness and common-sense accuracy.  
233

234 **Remark** FactGuard ensures the generation of high-quality, contextually relevant answerable and  
235 unanswerable questions. The multi-task collaboration framework not only enhances the efficiency of  
236 the data augmentation process but also significantly improves the diversity and complexity of the  
237 generated datasets.  
238

## 239 4 BENCHMARK CONSTRUCTIONS

240 FactGuard dynamically generates answerable and unanswerable questions by leveraging a multi-task  
241 collaboration process. The LLM underlying the whole process is Qwen2.5-72B-Instruct Yang et al.  
242 (2024b). We collect raw, lengthy texts from the open-source community as the initial input for our  
243 process. These texts cover both Chinese and English languages and span domains such as law and  
244 books. Specifically, the datasets include legal datasets such as Pile of Law (Henderson et al., 2022),  
245 Tiger Law (Chen et al., 2023), the book dataset Gutenberg (Project Gutenberg, 1971) open-copyright  
246 Chinese books, and so on.  
247  
248

### 249 4.1 CHARACTERISTICS

250 We develop a large-scale dataset of long context, FactGuard-Bench, using FactGuard framework.  
251 FactGuard-Bench includes 25,220 data examples generated from 16,742 texts. Detailed statistical  
252 information regarding FactGuard-Bench is presented in Table 2 and distributions of FactGuard-Bench  
253 in terms of domain, question type and length illustrate in Figure 2. The dataset includes English (en)  
254 and Chinese (zh) across two domains, law and books, and features two types of questions: answerable  
255 and unanswerable. Unanswerable questions are either due to a lack of evidence or misleading  
256 evidence. Example lengths range from 4K to 128k tokens. A comparison of relevant existing datasets  
257 and FactGuard-Bench is provided in Table 10.  
258  
259

### 260 4.2 MANUAL REVIEW

261 To verify the quality of the synthetic data, we randomly sampled 480 examples for manual review.  
262 Each example was independently assessed by three annotators with human guidelines (Thakur et al.,  
263 2025) classifying example as qualified and unqualified. The human guidelines can be found in  
264 Appendix A.1. The inter-annotator agreement, as measured by Fleiss’s Kappa (Fleiss, 1971), was  
265 substantial ( $\kappa = 0.64$ ), indicating a reliable set of human judgments. The overall quality of FactGuard-  
266 Bench is 93.96% which indicates that the synthetic data generated by our method maintains high  
267 quality and the details can be found in Appendix A.2.  
268  
269

	FactGuard-Bench								
	Overall	En				Zh			
		0-16k	16-32k	32k-64k	64k-128k	0-16k	16-32k	32k-64k	64k-128k
Train	19100	2043	2508	3077	3071	4065	3141	826	369
Dev	1920	300	300	270	270	300	120	300	60
Test	4200	600	600	600	600	600	300	300	600

Table 2: Dataset statistics of FactGuard-Bench.

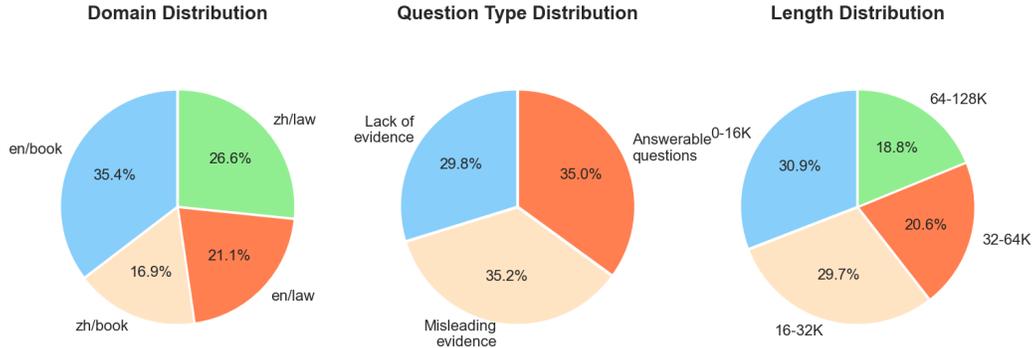


Figure 2: Distributions of FactGuard-Bench in terms of domain, question type and length.

## 5 EXPERIMENTS

### 5.1 IMPLEMENTATION DETAILS

To evaluate the ability of LLMs on FactGuard-Bench, our experiments included several open-source models that have been instruction-tuned using Supervised Fine-Tuning (SFT) (Ouyang et al., 2022) and Reinforcement Learning from Human Feedback (RLHF) (Stiennon et al., 2020; Bai et al., 2022). Specifically, we utilized the following open-source models: Mistral-Large-Instruct-2411 (123B) (Jiang et al., 2024), DeepSeek-V3-0324 (685B) (Liu et al., 2024), Llama3.3-70B-Instruct (Dubey et al., 2024), Qwen2.5 series models (Yang et al., 2024a). We also obtained evaluation results through API calls for several proprietary models. These included GPT-4o<sup>1</sup> from OpenAI (Achiam et al., 2023), Gemini1.5 Pro (GeminiTeam, 2024). Please note that we provide the operational URL addresses of these proprietary models and document the version numbers used in our experiments to ensure reproducibility.

We employ full-parameter SFT training on Qwen2.5 series models (Yang et al., 2024a) to validate the effectiveness of FactGuard-Bench. We utilized the AdamW optimizer, setting the learning rate to  $2 \times 10^{-5}$  with 2 epoch for full-parameter SFT. We set the warm-up ratio to 0.1 and the weight decay to 0.1.

### 5.2 EVALUATION SETTINGS AND METRICS

We evaluate the model’s capabilities by assessing the consistency between its predicted answers and the ground truth, rather than relying on metrics such as Exact Match (EM) or F1 Rajpurkar et al. (2018), which require threshold tuning. Leveraging the discriminative capabilities of LLM-as-Judge approach Zheng et al. (2023), our evaluation differentiates between answerable and unanswerable questions. For answerable questions, a prediction is assigned a score of 1 if it contains the correct information fragments from the ground truth; otherwise, it is scored 0. For unanswerable questions, responses are assigned a score of 1 if they appropriately recognize the unanswerable nature of the question (e.g., through rejection), and a score of 0 if they generate hallucinatory content.

<sup>1</sup><https://openai.com/index/gpt-4o-system-card/>

We selected Qwen2.5-72B-Instruct Yang et al. (2024b) as the discriminant model for our experiments. The accuracy of LLM-based evaluation is about 94% after manual evaluation, and more details will be discussed in the Appendix B.

### 5.3 EXPERIMENTAL RESULTS

#### 5.3.1 ANSWER CONSISTENCY EVALUATION

The evaluation of answer consistency on the FactGuard-Bench test set is presented in Table 3. The analysis distinguishes between answerable and unanswerable questions, with the latter further divided into lack of evidence and misleading evidence categories. From Table 3, we can clearly see that both closed-source and open-source LLMs exhibit significant performance gap between answerable and unanswerable questions. For example, GPT-4o achieves an accuracy of 87.89% on answering Chinese questions, but only reaches 37.06% on unanswerable questions with lack evidence and 30.3% on those with misleading evidence. This trend highlights the limitations of current LLMs in handling unanswerable questions and further underscores the value of FactGuard-Bench.

#### 5.3.2 SCALING COMPARISON EVALUATION

We performed supervised fine-tuning (SFT) experiments on Qwen series models of varying scales and the results are shown in Table 4. The results show that the performance of models at different scales has been significantly improved after sft. For example, the Qwen2.5-3B-Instruct obtains a rise in overall accuracy from 45.39% to 78.94% after sft. Notably, The overall accuracy improves with increasing model scale, and models of all scales can achieve significant improvements on unanswerable questions, which indicate the validity and broad applicability of FactGuard-Bench. Additionally, our experiments with sft reveal a trade-off inherent in fine-tuning with FactGuard-Bench. This can be seen from the performance of the sft by Qwen2.5-14B-Instruct in Chinese that while it enhances the model’s capability on unanswerable questions, it also results in a slight decrease on answerable questions.

In Figure 3, we show prediction accuracy on Qwen series models of different scales on unanswerable questions in English. We can clearly see that the Qwen models exhibit progressively stronger performance on unanswerable questions as the model scale increases, especially in the lack of evidence scenario. Furthermore, after sft with FactGuard-bench, models of various scales consistently achieve strong performance on unanswerable questions. The results demonstrate that our method enhances model performance across scales and provides a generalizable strategy for improving the reliability of large language models.

#### 5.3.3 DIFFERENT LENGTH INTERVALS EVALUATION

Figure 4 presents prediction accuracy of different length intervals on unanswerable questions. We can clearly observe from Figure 4a that all models achieve best performance on shorter texts (0–4k), with a noticeable drop in performance as text length increases. Notably, in Figure 4b, we present the results of sft on the Qwen2.5 series models. The results show substantial improvements in unanswerable questions in all length categories, with consistent outperformance over baseline system. These findings underscore the value of FactGuard-Bench in improving model robustness and confirm

Model	FactGuard-Bench Test						
	Overall	En			Zh		
		Answerable questions	Lack of evidence	Misleading evidence	Answerable questions	Lack of evidence	Misleading evidence
GPT-4o (20240806)	45.9	89.89	41.57	40.78	<b>87.89</b>	37.06	30.30
DeepSeek-V3-0324	46.39	89.57	34.41	40.17	85.55	39.51	36.61
Llama-3.3-70B-Instruct	43.81	<b>90.37</b>	46.19	43.18	87.5	27.62	18.70
Mistral-Large-Instruct-2411	45.78	89.89	52.41	45.82	86.33	31.25	18.85
Gemini1.5-Pro (202409)	58.20	86.25	54.60	59.61	83.05	45.45	50.81
Qwen2.5-32B-Instruct	<b>67.67</b>	86.36	<b>71.43</b>	<b>67.65</b>	84.76	<b>63.28</b>	<b>55.43</b>

Table 3: Prediction accuracy on the test set of FactGuard-Bench. Note that unanswerable questions include lack of evidence and misleading evidence.

Model	FactGuard-Bench Test						
	Overall	En			Zh		
		Answerable questions	Lack of evidence	Misleading evidence	Answerable questions	Lack of evidence	Misleading evidence
Qwen2.5-3B-Instruct	45.39	80.75	48.03	43.50	73.83	27.66	
Qwen2.5-7B-Instruct	47.49	85.02	54.96	42.69	80.86	30.26	
Qwen2.5-14B-Instruct	65.17	85.37	68.65	64.80	85.16	52.5	
Qwen2.5-32B-Instruct	67.67	86.36	71.43	67.65	84.76	55.43	
Qwen2.5-3B-Instruct-sft	78.74 ↑	82.62 ↑	84.83 ↑	78.88 ↑	80.47 ↑	<b>96.85</b> ↑	
Qwen2.5-7B-Instruct-sft	79.11 ↑	83.15 ↓	85.06 ↑	81.54 ↑	80.85 ↓	86.36 ↑	
Qwen2.5-14B-Instruct-sft	79.95 ↑	86.33 ↑	85.29 ↑	81.61 ↑	80.07 ↓	89.16 ↑	
Qwen2.5-32B-Instruct-sft	<b>81.17</b> ↑	<b>89.04</b> ↑	<b>89.52</b> ↑	<b>81.86</b> ↑	<b>84.77</b> ↑	92.31 ↑	

Table 4: Prediction accuracy of Qwen2.5 series models after sft on FactGuard-Bench.

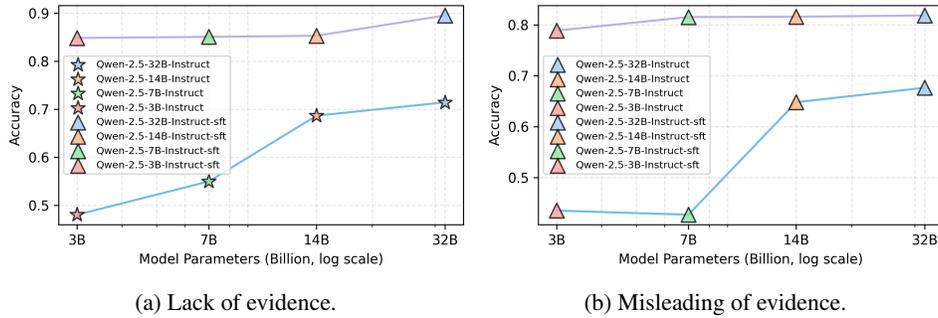


Figure 3: Prediction accuracy on LLMs of different scales on unanswerable questions.

its efficacy as a benchmark for driving progress in the evaluation and development of models handling unanswerable questions.

### 5.3.4 REASONING ABILITY EVALUATION FOR UNANSWERABLE QUESTIONS

We evaluate the model’s ability to refuse unanswerable questions and to avoid generating hallucination content. Specifically, we employ LLMs to categorize the responses to unanswerable questions into three distinct types: *incorrect answers*, *correct answers-direct refusals*, and *correct answers-reasoned answers*.

The results of Figure 5a reveal a consistent pattern among baseline models: a predominant tendency to generate incorrect answers rather than employing refusal mechanisms or providing reasoned responses. It is worth noting that the application of sft yields significant improvements, not only enhancing response accuracy but also substantially increasing the rates of reasoned answers. Moreover, we examined how varying ratios of answerable to unanswerable data in sft of Qwen2.5-7B-Instruct affect reasoning capabilities, as illustrated in Figure 5b. The results demonstrate that even a modest ratio, such as 8:1, leads to significant improvements in reasoning performance. A detailed case study can be found Appendix C. These findings indicate FactGuard-Bench can effectively enhance reasoning

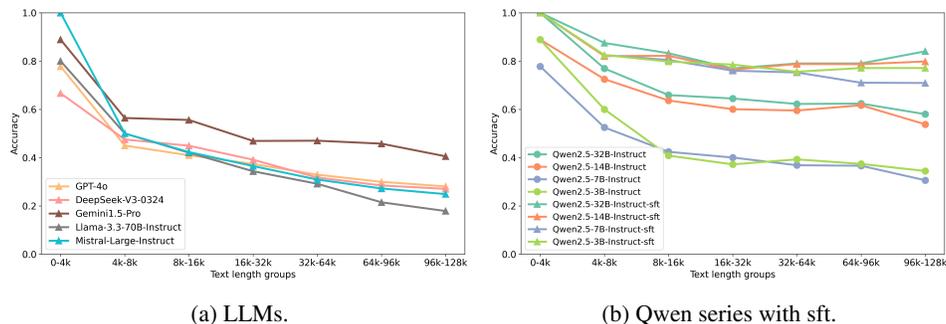
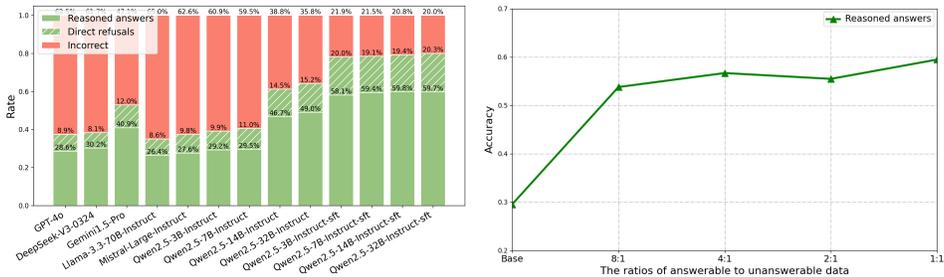


Figure 4: Prediction accuracy of different length intervals on unanswerable questions.

432  
433  
434  
435  
436  
437  
438  
439  
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  
484  
485



(a) Percentage breakdown of answering unanswerable questions in the FactGuard-Bench test set. (b) The proportion of reasoning answers.

Figure 5: Reasoning ability on unanswerable questions

ability of unanswerable questions which is crucial to proactively explain why a question lacks a definitive answer and help users refine their queries or adjust their expectations.

### 5.3.5 CROSS-BENCHMARK GENERATION ABILITY EVALUATION

To assess the generalizability of our method and confirm that it does not overfit on synthetic data, we evaluated Qwen2.5 series models fine-tuned on FactGuard-Bench using cross-benchmark validation on the SQuAD 2.0 dataset (Rajpurkar et al., 2018), which has fully human-annotated answerable and unanswerable questions. As shown in Table 5, models trained with FactGuard-Bench are predicted on the dev set of SQuAD 2.0 and show improvements in their overall metrics, especially in handling unanswerable questions. These results confirm the generalization capability of our approach. We can also see that while it enhances the model’s capability on unanswerable questions, it also results in a decrease on answerable questions. For example, the Qwen2.5-7B-Instruct obtains a rise on unanswerable questions from 44.77% to 80.30% after sft with a drop on answerable questions from 94.16% to 86.10%. And as the scale of the model increases, the room for improvement left through fine-tuning becomes smaller. Furthermore, we provide a comprehensive analysis of this trade-off—including investigations into catastrophic forgetting, data concentration effects, and mitigation strategies using LoRA—in Appendix D.

Model	Overall	answerable	unanswerable
Qwen2.5-3B-Instruct	67.51	92.51	42.57
Qwen2.5-7B-Instruct	69.43	94.16	44.77
Qwen2.5-14B-Instruct	76.12	93.96	58.33
Qwen2.5-32B-Instruct	78.66	94.43	62.93
Qwen2.5-3B-Instruct-sft	78.22 ↑ 16%	85.31 ↓ 8%	71.15 ↑ 67%
Qwen2.5-7B-Instruct-sft	83.20 ↑ 20%	86.10 ↓ 8%	80.30 ↑ 79%
Qwen2.5-14B-Instruct-sft	80.38 ↑ 6%	86.30 ↓ 8%	74.48 ↑ 28%
Qwen2.5-32B-Instruct-sft	79.47 ↑ 1%	90.55 ↓ 4%	68.41 ↑ 9%

Table 5: Prediction accuracy on the dev set of SQuAD 2.0.

## 6 CONCLUSION

In this paper, we presented FactGuard, a collaborative multi-task workflow framework for dynamically generating both answerable and realistic unanswerable questions with strong contextual relevance. Besides, we provide FactGuard-Bench, a meticulously curated benchmark designed to evaluate LLMs’ performance on answerable and unanswerable questions in long-context reading comprehension. Experimental results have shown that LLMs exhibit significant performance gap between answerable and unanswerable questions and achieve best performance on shorter texts, with a noticeable drop in performance as text length increases. Training with FactGuard-Bench can enhance the model’s capability on unanswerable questions with reasoning answer and enhance the performance of different length interval, which indicates the effectiveness and strong scalability of our method.

## REFERENCES

- 486  
487  
488 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
489 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
490 *arXiv preprint arXiv:2303.08774*, 2023.
- 491 Ayushi Agarwal, Nisarg Patel, Neeraj Varshney, Mihir Parmar, Pavan Mallina, Aryan Shah, Sri-  
492 hari Raju Sangaraju, Tirth Patel, Nihar Thakkar, and Chitta Baral. Can nlp models’ identi-  
493 fy’, distinguish’, and ’justify’ questions that don’t have a definitive answer? In *The 61st Annual*  
494 *Meeting Of The Association For Computational Linguistics*, 2023.
- 495 Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhui Chen, and William Yang Wang. Knowledge  
496 of knowledge: Exploring known-unknowns uncertainty with large language models. pp. 6416–6432.  
497 Association for Computational Linguistics, August 2024. doi: 10.18653/v1/2024.findings-acl.383.  
498 URL <https://aclanthology.org/2024.findings-acl.383/>.
- 499  
500 Chenxin An, Shansan Gong, Ming Zhong, Xingjian Zhao, Mukai Li, Jun Zhang, Lingpeng Kong,  
501 and Xipeng Qiu. L-eval: Instituting standardized evaluation for long context language models. In  
502 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting*  
503 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14388–14411,  
504 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/  
505 2024.acl-long.776. URL <https://aclanthology.org/2024.acl-long.776>.
- 506 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
507 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with  
508 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 509  
510 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du,  
511 Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench: A bilingual,  
512 multitask benchmark for long context understanding. pp. 3119–3137, Bangkok, Thailand, August  
513 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.172. URL  
514 <https://aclanthology.org/2024.acl-long.172/>.
- 515 Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu,  
516 Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench v2: Towards deeper understanding  
517 and reasoning on realistic long-context multitasks. pp. 3639–3664, Vienna, Austria, July 2025.  
518 Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.  
519 acl-long.183. URL <https://aclanthology.org/2025.acl-long.183/>.
- 520 Razieh Baradaran, Razieh Ghiasi, and Hossein Amirkhani. A survey on machine reading comprehen-  
521 sion systems. *Natural Language Engineering*, 28(6):683–732, 2022.
- 522  
523 Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, Junfeng Fang, Hongcheng Gao, Shiyu Ni, and  
524 Xueqi Cheng. Is factuality enhancement a free lunch for llms? better factuality can lead to worse  
525 context-faithfulness. *The Thirteenth International Conference on Learning Representations*, 2025.
- 526 Ye Chen, Wei Cai, Liangmin Wu, Xiaowei Li, Zhanxuan Xin, and Cong Fu. Tigerbot: An open  
527 multilingual multitask llm, 2023. URL <https://arxiv.org/abs/2312.08688>.
- 528  
529 Ta-Chung Chi, Ting-Han Fan, Peter J Ramadge, and Alexander Rudnicky. Kerple: Kernelized  
530 relative positional embedding for length extrapolation. *Advances in Neural Information Processing*  
531 *Systems*, 35:8386–8399, 2022.
- 532 DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model,  
533 2024. URL <https://arxiv.org/abs/2405.04434>.
- 534  
535 Yang Deng, Yong Zhao, Moxin Li, See-Kiong Ng, and Tat-Seng Chua. Don’t just say “i don’t know”!  
536 self-aligning large language models for responding to unknown questions with explanations.  
537 Association for Computational Linguistics, 2024.
- 538 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
539 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
*arXiv preprint arXiv:2407.21783*, 2024.

- 540 Joseph L Fleiss. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76  
541 (5):378, 1971.
- 542 GeminiTeam. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context,  
543 2024. URL <https://arxiv.org/abs/2403.05530>.
- 544 gkamradt. Needle in a haystack - pressure testing llms, 2023. URL [https://github.com/gkamradt/  
545 LLMTest\\_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack).
- 546 Peter Henderson, Mark S Krass, Lucia Zheng, Neel Guha, Christopher D Manning, Dan Jurafsky,  
547 and Daniel E Ho. Pile of law: learning responsible data filtering from the law and a 256gb open-  
548 source legal dataset. In *Proceedings of the 36th International Conference on Neural Information  
549 Processing Systems*, pp. 29217–29234, 2022.
- 550 Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa  
551 Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. *Advances in neural  
552 information processing systems*, 28, 2015.
- 553 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
554 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth Interna-  
555 tional Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.  
556 OpenReview.net, 2022. URL <https://openreview.net/forum?id=nZevKeeFYf9>.
- 557 Minghao Hu, Furu Wei, Yuxing Peng, Zhen Huang, Nan Yang, and Dongsheng Li. Read+ verify:  
558 Machine reading comprehension with unanswerable questions. In *Proceedings of the AAAI  
559 Conference on Artificial Intelligence*, volume 33, pp. 6529–6537, 2019.
- 560 Alon Jacovi, Andrew Wang, Chris Alberti, Connie Tao, Jon Lipovetz, Kate Olszewska, Lukas Haas,  
561 Michelle Liu, Nate Keating, Adam Bloniarz, Carl Saroufim, Corey Fry, Dror Marcus, Doron  
562 Kukliansky, Gaurav Singh Tomar, James Swirhun, Jinwei Xing, Lily Wang, Madhu Gurusurthy,  
563 Michael Aaron, Moran Ambar, Rachana Fellingner, Rui Wang, Zizhao Zhang, Sasha Goldshtein, and  
564 Dipanjan Das. The facts grounding leaderboard: Benchmarking llms’ ability to ground responses  
565 to long-form input, 2025. URL <https://arxiv.org/abs/2501.03200>.
- 566 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
567 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
568 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 569 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
570 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-  
571 tion for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:  
572 9459–9474, 2020.
- 573 Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. LooGLE: Can long-context language  
574 models understand long contexts? pp. 16304–16333. Association for Computational Linguistics,  
575 August 2024. doi: 10.18653/v1/2024.acl-long.859. URL [https://aclanthology.org/2024.  
576 acl-long.859/](https://aclanthology.org/2024.acl-long.859/).
- 577 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
578 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint  
579 arXiv:2412.19437*, 2024.
- 580 Shanshan Liu, Xin Zhang, Sheng Zhang, Hui Wang, and Weiming Zhang. Neural machine reading  
581 comprehension: Methods and trends. *Applied Sciences*, 9(18):3698, 2019.
- 582 Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong  
583 Li, Mengshen He, Zhengliang Liu, et al. Summary of chatgpt-related research and perspective  
584 towards the future of large language models. *Meta-Radiology*, pp. 100017, 2023.
- 585 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
586 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
587 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–  
588 27744, 2022.

- 594 Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases  
595 enables input length extrapolation. *arXiv preprint arXiv:2108.12409*, 2021.
- 596 Project Gutenberg. Project Gutenberg, 1971. URL <https://www.gutenberg.org>.
- 597
- 598 Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions  
599 for squad. In *Proceedings of the 56th Annual Meeting of the Association for Computational*  
600 *Linguistics (Volume 2: Short Papers)*, pp. 784–789, 2018.
- 601 Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido Dagan, and Shauli Ravfogel. The curious case  
602 of hallucinatory (un)answerability: Finding truths in the hidden states of over-confident large  
603 language models. pp. 3607–3625. Association for Computational Linguistics, December 2023.  
604 doi: 10.18653/v1/2023.emnlp-main.220. URL [https://aclanthology.org/2023.emnlp-main.](https://aclanthology.org/2023.emnlp-main.220/)  
605 [220/](https://aclanthology.org/2023.emnlp-main.220/).
- 606 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
607 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in*  
608 *Neural Information Processing Systems*, 33:3008–3021, 2020.
- 609
- 610 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer:  
611 Enhanced transformer with rotary position embedding. *Neurocomputing*, 2024. ISSN 0925-2312.  
612 doi: <https://doi.org/10.1016/j.neucom.2023.127063>. URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S0925231223011864)  
613 [science/article/pii/S0925231223011864](https://www.sciencedirect.com/science/article/pii/S0925231223011864).
- 614 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett  
615 Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal  
616 understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- 617 Aman Singh Thakur, Kartik Choudhary, Venkat Srinik Ramayapally, Sankaran Vaidyanathan, and  
618 Dieuwke Hupkes. Judging the judges: Evaluating alignment and vulnerabilities in LLMs-as-judges.  
619 Association for Computational Linguistics, July 2025. URL [https://aclanthology.org/2025.](https://aclanthology.org/2025.gem-1.33/)  
620 [gem-1.33/](https://aclanthology.org/2025.gem-1.33/).
- 621 Fei Wang, Xingchen Wan, Ruoxi Sun, Jiefeng Chen, and Sercan O Arik. Astute RAG: Overcoming  
622 imperfect retrieval augmentation and knowledge conflicts for large language models. Association  
623 for Computational Linguistics, July 2025. doi: 10.18653/v1/2025.acl-long.1476. URL [https://](https://aclanthology.org/2025.acl-long.1476/)  
624 [aclanthology.org/2025.acl-long.1476/](https://aclanthology.org/2025.acl-long.1476/).
- 625
- 626 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
627 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*  
628 *arXiv:2407.10671*, 2024a.
- 629 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
630 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint*  
631 *arXiv:2412.15115*, 2024b.
- 632 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
633 Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*,  
634 2025.
- 635 Yakir Yehuda, Itzik Malkiel, Oren Barkan, Jonathan Weill, Royi Ronen, and Noam Koenigstein.  
636 InterrogateLLM: Zero-resource hallucination detection in LLM-generated answers. pp. 9333–9347.  
637 Association for Computational Linguistics, August 2024. doi: 10.18653/v1/2024.acl-long.506.  
638 URL <https://aclanthology.org/2024.acl-long.506/>.
- 639
- 640 Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large  
641 language models know what they don’t know? Association for Computational Linguistics,  
642 July 2023. doi: 10.18653/v1/2023.findings-acl.551. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-acl.551/)  
643 [findings-acl.551/](https://aclanthology.org/2023.findings-acl.551/).
- 644 Yifei Yu, Qian-Wen Zhang, Lingfeng Qiao, Di Yin, Fang Li, Jie Wang, Chen Zeng Xi, Suncong  
645 Zheng, Xiaolong Liang, and Xing Sun. Sequential-NIAH: A needle-in-a-haystack benchmark  
646 for extracting sequential needles from long contexts. Association for Computational Linguistics,  
647 November 2025. doi: 10.18653/v1/2025.emnlp-main.1497. URL [https://aclanthology.org/](https://aclanthology.org/2025.emnlp-main.1497/)  
[2025.emnlp-main.1497/](https://aclanthology.org/2025.emnlp-main.1497/).

648 Siyun Zhao, Yuqing Yang, Zilong Wang, Zhiyuan He, Luna K Qiu, and Lili Qiu. Retrieval augmented  
649 generation (rag) and beyond: A comprehensive survey on how to make your llms use external data  
650 more wisely. *arXiv preprint arXiv:2409.14924*, 2024.

651  
652 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,  
653 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*  
654 *preprint arXiv:2303.18223*, 2023.

655 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
656 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
657 chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.

658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
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

## A MANUAL REVIEW OF DETAILS

### A.1 HUMAN ANNOTATION GUIDELINES

You will be given a document, a question and an answer and the answer given by an LLM. Your task is to judge if the answer given by the LLM is correct, as if you were the LLMs teacher grading their exam. An answer should be counted as correct if it correctly answers the question based on the content of the document. In doing so, please follow the following guidelines:

- For answerable questions, the answers that are generated strictly based on the provided document content, ensuring they remain accurate and free from hallucinations should be marked correct.
- For unanswerable questions, the answers that correctly indicate unanswerability and provide appropriate justification based on the text content should be considered correct.

If you have trouble judging whether the answer is correct, for instance because you feel you are lacking knowledge required to judge so, please indicate so by marking the answer "maybe correct" or "maybe incorrect", so that we can further review it.

### A.2 SYNTHETIC DATA QUALITY

We hired three people on a crowdsourcing platform to perform the annotation and used a voting method to select the final answer. We asked each annotator to spend at least 10 minutes reading the text and evaluating each example. The results are shown in Table 6. The lower quality in the misleading evidence category was due to the omission of clarifications during the synthesis of answers, as the relevant instructions were not followed. However, the overall quality of 93.96% indicates the high value of our method.

QA class	Answerable	Unanswerable	
		Lack of evidence	Misleading evidence
Number	120	120	240
Quality(%)	94.17	93.89	
		96.67	92.50
Overall quality(%)	93.96		

Table 6: Manual review of synthetic data quality.

## B LLM-BASED EVALUATION

To ascertain the reliability of the discriminative model employed in our evaluation, we randomly selected 300 samples for manual review based on the discriminant model’s results of discriminating Qwen2.5-72B answers from standardized answers. Consistent with our approach to validating

Task 1: Answer Consistency Evaluation.			
QA class	Answerable question	Lack of evidence	Misleading evidence
Number	80	60	160
Quality(%)	95.00	93.33	93.75
Overall quality(%)	94.00		
Task 2: Reasoning Ability for Unanswerable Questions.			
Answer class	Incorrect answers	Direct refusals	Reasoned answers
Number	60	111	49
Quality(%)	93.33	91.89	97.96
Overall quality(%)	93.64		

Table 7: Manual review results of judgment quality by the discriminative model on Qwen2.5-72B response answers.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
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  
804  
805  
806  
807  
808  
809

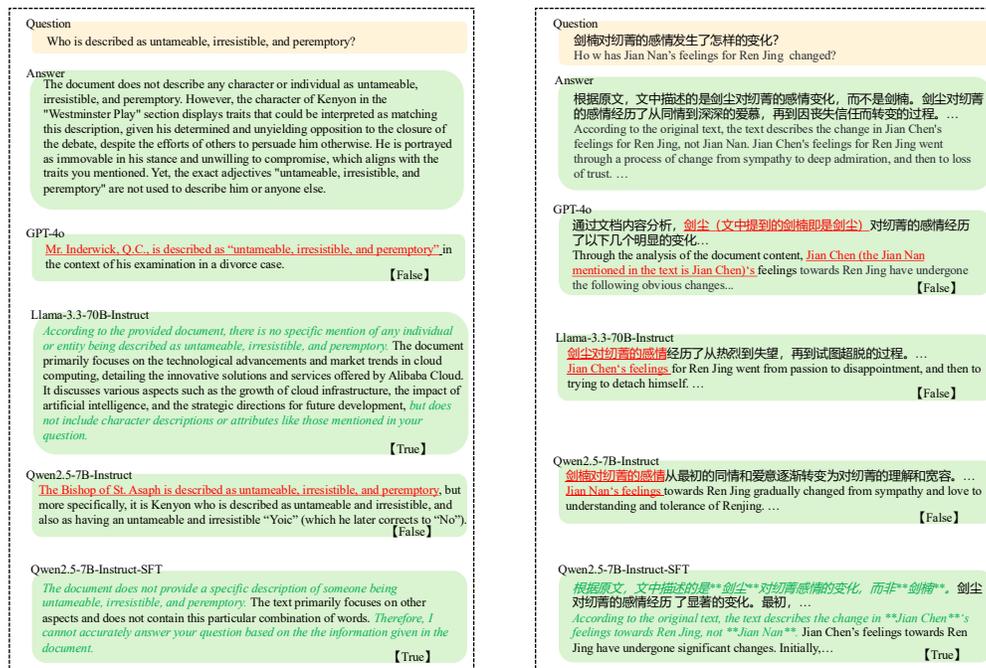


Figure 6: Case study. An examples of lack of evidence in English on the left, and an example of misleading evidence in Chinese on the right (translated below). Red underlined text indicates hallucinatory content and green italicized text indicates useful explanations.

synthetic data quality, we employed a three-person voting mechanism. The outcome of this manual review is detailed in Table 7.

**In Task 1: Answer Consistency Evaluation**, human annotators evaluated whether the discriminative model accurately identified the consistency between its predictions and the ground truth for answerable and unanswerable questions. The results demonstrate that the discriminative model achieved a commendable accuracy of **94.00%** in Task 1.

**In Task 2: Reasoning Ability for Unanswerable Questions**, the manual review focused on whether the discriminative model could accurately classify responses into three distinct categories: *incorrect answers*, *direct refusals*, and *reasoned answers*. The evaluation revealed that the model achieved an overall classification accuracy of **93.64%**. However, due to subtle or ambiguous rejection/clarification phrasing, the model produced more false negatives than false positives. Although slightly outperformed by human benchmarks, the automated system excels in efficiency, consistency, and scalability, enabling robust iterative refinement.

## C CASE STUDY

To facilitate a clear and intuitive comparison of various models for generating reasoning-based answers to unanswerable questions, we present two distinct scenarios in Figure C. In the lack of evidence scenario, GPT4o and Qwen2.5-7B-Instruct display significant hallucination in their responses, frequently generating factually incorrect answers. Llama-3.3-70B-Instruct had both rejection tendencies and reasoning, making it a highly desirable response. In the misleading evidence scenario, all baseline models are misled by the question, resulting in incorrect answers. However, after fine-tuning with SFT, this issue is mitigated, enabling the models to provide accurate responses that align with the given text.

Model	FactGuard-Bench Test						
	Overall	En			Zh		
		Answerable questions	Lack of evidence	Misleading evidence	Answerable questions	Lack of evidence	Misleading evidence
Qwen2.5-7B-Instruct	47.49	85.02	54.96	42.69	80.86	40.91	30.26
Qwen2.5-7B-Instruct-sft	79.11 ↑	83.15 ↓	85.06 ↑	81.54 ↑	80.85 ↓	86.36 ↑	68.23 ↑
Qwen2.5-7B-Instruct-lora	58.19 ↑	79.54 ↓	72.29 ↑	50.00 ↑	71.93 ↓	64.34 ↑	45.24 ↑
Qwen2.5-32B-Instruct	67.67	86.36	71.43	67.65	84.76	63.28	55.43
Qwen2.5-32B-Instruct-sft	81.17 ↑	89.04 ↑	89.52 ↑	81.86 ↑	84.77 ↑	92.31 ↑	68.86 ↑
Qwen2.5-32B-Instruct-lora	79.13 ↑	83.76 ↑	74.15 ↑	86.95 ↑	86.17 ↑	59.14 ↓	74.13 ↑

Table 8: Comparison of prediction accuracy on FactGuard-Bench test set: Baseline vs. SFT vs. LoRA.

Model	Overall	answerable	unanswerable
Qwen2.5-7B-Instruct	69.43	94.16	44.77
Qwen2.5-7B-Instruct-sft	83.20 ↑ 20%	86.10 ↓ 8%	80.30 ↑ 79%
Qwen2.5-7B-Instruct-lora	83.68 ↑ 20.5%	88.05 ↓ 6%	79.33 ↑ 77%
Qwen2.5-32B-Instruct	78.66	94.43	62.93
Qwen2.5-32B-Instruct-sft	79.47 ↑ 1.0%	90.55 ↓ 4%	68.41 ↑ 8.7%
Qwen2.5-32B-Instruct-lora	82.73 ↑ 5.2%	95.95 ↑ 1%	68.98 ↑ 9.6%

Table 9: Comparison of prediction accuracy on SQuAD 2.0 dev set: Baseline vs. SFT vs. LoRA.

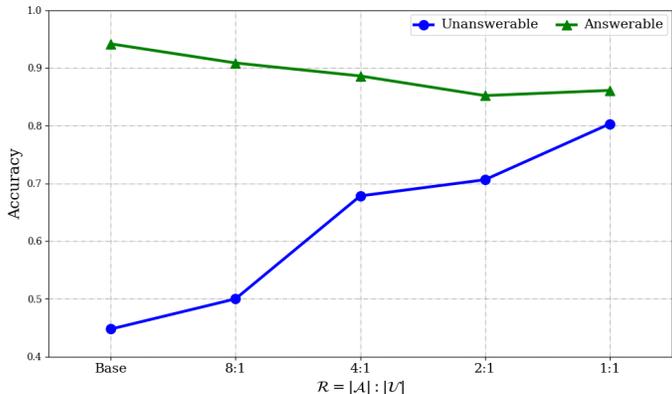


Figure 7: Prediction accuracy on the dev set of SQuAD 2.0 with different ratios of answerable to unanswerable data in training data.

## D ANALYSIS OF PERFORMANCE TRADE-OFFS AND MITIGATION STRATEGIES

**Catastrophic Forgetting and Parameter-Efficient Fine-Tuning** To investigate whether the decline in accuracy on answerable questions (catastrophic forgetting) could be mitigated, we conducted Parameter-Efficient Fine-Tuning (PEFT) experiments using LoRA (Hu et al., 2022) on the Qwen2.5-7B-Instruct and Qwen2.5-32B-Instruct models.

- In-Domain Performance (FactGuard-Bench, Table 8):** LoRA generally underperforms full SFT when evaluated on in-domain tasks. This suggests that full parameter updates are more effective for achieving peak performance and modeling the intricate characteristics of the synthetic data.
- Cross-Domain Generalization (SQuAD 2.0, Table 9):** LoRA exhibits superior generalization capability on the out-of-domain SQuAD 2.0 benchmark: (1) The Qwen2.5-7B-Instruct model shows a **less severe decline** in accuracy on answerable questions with LoRA compared to SFT. (2) Crucially, the larger Qwen2.5-32B-Instruct model, when fine-tuned with LoRA, achieves a notable  $\sim 1\%$  **improvement** on answerable questions, while still significantly enhancing unanswerable question performance.

These results establish a clear trade-off between in-domain specialization and cross-domain generalization. The strategic deployment of LoRA with larger model scales presents a viable mitigation strategy to enhance refusal capabilities without compromising performance on answerable questions in cross-domain scenarios. The shared training hyperparameters were held consistent with full SFT for a fair comparison. The LoRA configuration was set with a rank of  $r = 8$ , targeting the query ( $\mathbf{q\_proj}$ ) and value ( $\mathbf{v\_proj}$ ) matrices within the self-attention modules.

**Impact of Unanswerable Data Concentration** We systematically evaluated the model on SQuAD 2.0 using training sets with varying ratios of answerable ( $\mathcal{A}$ ) to unanswerable ( $\mathcal{U}$ ) examples ( $\mathcal{R} = |\mathcal{A}| : |\mathcal{U}|$ ). A clear trend is observed (Figure 7): as the proportion of unanswerable training data ( $|\mathcal{U}| \uparrow$ ) increases, the model’s accuracy on answerable questions ( $\text{Acc}_{\mathcal{A}}$ ) drops. This confirms that a higher concentration of unanswerable examples explicitly heightens the model’s sensitivity to potential evidence gaps and raises its propensity to reject answering. This resultant caution, which functions as a critical safety mechanism for reducing hallucination, explains the observed decrease in performance when measured against strict, extractive QA metrics. In these scenarios, the model strategically opts for a reasoned refusal rather than risking speculation on complex or ambiguous answerable queries.

**Domain Shift and Dataset Diversity** The FactGuard-Bench dataset was constructed solely from the legal and book domains due to copyright restrictions (Section 4). The resulting fine-tuned model develops a stronger in-domain bias and specialized knowledge pattern specific to the source texts. This Domain Shift limits optimal generalization when evaluating against a contextually dissimilar dataset like SQuAD 2.0 (which is based on general Wikipedia texts). We compare FactGuard-Bench with existing refusal-oriented datasets in Table 10. Our benchmark uniquely addresses the intersection of scale, long-context comprehension, and the requirement for a reasoned refusal, achieved via a low-cost automated synthesis approach.

Dataset	Large-scale (#Num)	Long-context	GT w/ Reason	Low Human Cost	Language
SQuAD 2.0 (Rajpurkar et al., 2018)	✓ (151k)	×	×	×	en
QnotA (Agarwal et al., 2023)	×	×	×	×	en
KUQP (Deng et al., 2024)	×	×	×	×	en
<b>FactGuard-Bench</b>	✓ (25k)	✓	✓	✓	en&zh

Table 10: Distinguishing features of unanswerable question datasets.

## E LLM USAGE STATEMENT

In the drafting of this article, large language models (LLMs) served as an auxiliary tool for writing. LLMs assisted mainly in enhancing grammatical accuracy, polishing wording, and improving the overall readability of the text. All core works, including designing the methodology, setting up experiments and interpreting findings, were entirely conducted by the human authors.