ABG-SCIQA: A DATASET FOR UNDERSTANDING AND RESOLVING AMBIGUITY IN SCIENTIFIC QUESTIONS

Anonymous authors

Paper under double-blind review

Abstract

Asking ambiguous questions is a natural aspect of human communication, making it essential for Large Language Models (LLMs) to effectively recognize and address ambiguities. However, there is a lack of a comprehensive analysis of how well LLMs detect and solve ambiguities. Moreover, while several datasets on ambiguity exist, the absence of explicit explanations of ambiguity and annotations of ambiguity types limits the comprehensive evaluation. To address this issue, we introduce Abg-SciOA, a dataset designed to evaluate and help LLMs detect ambiguities and generate appropriate clarification questions using challenge questions in the area of social and nature science. Abg-SciQA encompasses four tasks: Ambiguity Detection, Ambiguity Type Classification, Clarification Question Generation, and Clarification-Based Question Answering, where each task has corresponding annotations. We evaluate the dataset using both closed-source and opensource LLMs and fine-tune open-source models. Our experiments show that the most state-of-the-art LLMs still encounter difficulties in resolving ambiguity in natural questions, and fine-tuning on Abg-SciQA can significantly enhance their capabilities to understand and address ambiguities. Notably, in the Ambiguity Type Classification task, the F1 score of Llama2-13b improves significantly from 16.6% to 79.1%. On the other hand, Abg-SciQA remains a challenging benchmark for LLMs, revealing ample room for model improvement. Our dataset can be found here¹.

029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

1 INTRODUCTION

Large Language Models (LLMs) have become widely used in various applications, including conversational systems (Achiam et al., 2023), code generation (Du et al., 2024), and optimization (Yu et al., 2023). However, LLMs often face challenges when dealing with ambiguous questions—questions that can have multiple interpretations or unclear meanings. For the first example in Fig. 2, the question "What?" could refer to "What did he do for the children?", or simply asking for a repeat for the previous answer, depending on the context. Such ambiguity makes it difficult for LLMs to provide accurate answers, as they may exhibit overconfidence in their responses (Xiong et al., 2023). Given that ambiguous questions are common in natural human communication (Clark & Brennan, 1991), addressing this issue is crucial for improving LLM performance.

To address ambiguity in natural language processing, researchers have focused on generating clari-043 fication questions as a key strategy. Language models are often used to automatically generate these 044 questions to resolve ambiguities (Zamani et al., 2020; Deng et al., 2023b), while other approaches rely on predefined clarification questions (Eberhart & McMillan, 2022; Aliannejadi et al., 2019). 046 The success of these methods is heavily dependent on the quality of the datasets used. High-quality 047 datasets are crucial not only for producing accurate clarification questions, but also for enhancing 048 the overall ability of LLMs to manage ambiguous queries. Several datasets have been developed with this goal in mind. For instance, Abg-CoQA (Guo et al., 2021) is an extension of CoQA (Reddy et al., 2019) that includes ambiguous questions and related clarifications. Similarly, AmbigQA (Min 051 et al., 2020) is derived from NQ-Open (Lee et al., 2019). Other datasets, such as the one proposed by Rao & Daumé III (2018), use StackExchange as source data. 052

¹https://anonymous.4open.science/r/Abg-Sci-DF10/README.md



Figure 1: An example of a data sample in Abg-SciQA. Each sample in Abg-SciQA includes a story and a corresponding question, covering four tasks: 1) Ambiguity Detection, 2) Ambiguity Type Classification, 3) Clarification Question Generation, and 4) Clarification-Based Question Answering. Unlike previous datasets, Abg-SciQA features an additional task for classifying types of ambiguity, enabling a more comprehensive analysis.

076 However, existing benchmarks have the follow-077 ing limitations. First, question-answers in these source datasets, such as CoQA (Reddy et al., 2019), are publicly available and thus may be 079 part of the LLM pertaining data mixture. As a result, evaluating on those benchmarks may not 081 reveal the models' real capabilities in addressing the ambiguity. Second, ambiguity annota-083 tions in dialogue-based datasets are sometimes 084 questionable. Given the characteristics of con-085 tinuity in dialogues, a question is usually considered as non-ambiguous even if the explicit 087 reference is missing. Taking the first example in Fig. 2, it is obvious that "What" refers to "What did he do for the children?" considering 089 the dialogue continuity, thus non-ambiguous. 090 Similarly, the second example should also be 091 labeled as non-ambiguous since the book right 092 after "drawing of her mother" is in red. Third, many of these datasets lack detailed annotations 094 for different types of ambiguity, limiting their effectiveness in broader evaluations. 096

075

To address these limitations, we introduce a new dataset on ambiguity, Abg-SciQA, which leverages the capabilities of LLMs and incorporates articles from various natural and social



Figure 2: An example of quality issues in previous datasets. Abg-CoQA (Guo et al., 2021) follows a conversational format. The first question considers dialogue continuity, resulting in an unambiguous query, while the second question lacks this consideration, leading to ambiguity. This inconsistency may confuse both humans and LLMs regarding the dataset's standards.

100 science domains. Resolving ambiguity in these fields is crucial to ensuring precision and accuracy 101 in science, enabling clear communication and collaboration across disciplines. It also improves 102 decision-making and ethical standards, particularly in social sciences, where ambiguous questions 103 can lead to flawed or harmful outcomes. Ultimately, clarity in questions advances the accumulation 104 of reliable knowledge in both the natural and social sciences, driving progress and understanding. 105 To avoid overlapping with pretraining data, ambiguous questions in Abg-SciQA are automatically generated by LLM. Then, an auxiliary LLM is then employed to assess the quality of the generated 106 samples, with human evaluators also involved in the evaluation process. Finally, Abg-SciQA classi-107 fies each ambiguity question into four distinct types of ambiguity, in addition to the standard detec-

Dataset	Data Size	# of Entries	Abg Rate	Ambiguity Detection	Type Classification	Clarification Generation	Clarification-Base QA
AmbigQA (Min et al., 2020)	64.0M	14,042	51.1%	1	X	X	X
Abg-CoQA (Guo et al., 2021)	21.1M	8,615	11.5%	1	X	1	1
ASQA (Stelmakh et al., 2022)	14.0M	6,316	45%	1	X	X	X
CAMBIGNO (Lee et al., 2023)	27.8M	5,653	100%	1	X	1	1
Abg-SciQA (Ours)	52.3M	13,729	73.3%	1	1	1	1

Table 1: Comparisons of different datasets on ambiguity. The "# of Entries" means the total number
of questions in the whole dataset. The "Abg Rate" represents the percentage of ambiguous questions
relative to the total number of questions. Abg-SciQA comprises four different tasks. Compared to
other datasets on ambiguity, our dataset covers the widest range of tasks with decent numbers of
entries and Abg Rate.

Ambiguity Type	Definition
Lexical	This occurs when a word has multiple meanings or multiple interpretations.
Syntactic	This arises from the structure of a sentence. A sentence can have multiple interpretations depending on how it's parsed.
Incomplete	This occurs when a statement or question lacks essential contextual information—such as location, time, event, or people—resulting in multiple possible interpretations.
Contextual	This type of ambiguity occurs when a question is clear and unambiguous in its wording, but it contains two possible answers due to differing contexts, interpretations, or sources.

Table 2: The definition of four types of ambiguity: Lexical, Syntactic, Incomplete, and Contextual Ambiguity.

tion and clarification tasks. This type-based annotation enables a more compressive analysis. Fig. 1
shows a data instance in Abg-SciQA. In addition, each instance in Abg-SciQA consists of a unique
question for the story, which avoids confusion brought by dialogues. We compare Abg-SciQA with
other datasets in the ambiguity area in Table 1. Specifically, Abg-SciQA includes 1,353 passages
and 13,729 questions, of which 10,202 are ambiguous, which is the largest dataset in the area. In this
paper, we evaluate Abg-SciQA on both closed-source LLMs and open-source LLMs. We outline our
contributions as follows:

- We introduce Abg-SciQA, a dataset that includes challenging ambiguous questions from diverse scientific fields, complete with annotations for different types of ambiguity.
- To the best of our knowledge, we are the first to establish benchmarks for solving ambiguous questions using both closed-source LLMs, such as GPT-o1, and open-source LLMs, such as Llama2 (Touvron et al., 2023).
- Our comprehensive experiments demonstrate that fine-tuning LLMs with Abg-SciQA can significantly improve their ability to handle ambiguous questions.
- 144 145 146

136

137

138

139

140

141

142

143

2 RELATED WORK

147 **Dealing with Ambiguity** Detecting and resolving ambiguous questions is crucial for modern language models and dialogue systems (Deng et al., 2023a). Significant efforts have been made in 148 this area. For instance, the skill of requesting clarification in dialogue has been extensively studied 149 by Purver et al. (2003), Schlangen (2004), and Stoyanchev et al. (2014). This task generally in-150 volves two subtasks (Aliannejadi et al., 2021): 1) detecting ambiguous questions and 2) generating 151 clarification questions. Typically, detecting ambiguous questions is treated as a binary task, where 152 the model determines whether a given question is ambiguous (Guo et al., 2021). Zhang & Choi 153 (2023) addresses this problem using uncertainty estimation. Various approaches have been pro-154 posed to generate clarification questions. Khalid et al. (2020) and Rodríguez & Schlangen (2004) 155 integrate cognitive modeling and discourse theories with reinforcement learning to select the most 156 effective clarification requests from a set of predefined strategies. Similarly, Eberhart & McMillan 157 (2022) and Aliannejadi et al. (2019) sample predefined clarification questions based on the ambigu-158 ous question. Another line of work focuses on generating clarification questions using language models (Guo et al., 2021; Zamani et al., 2020). Deng et al. (2023b) propose a method for generating 159 clarification questions through spatial reasoning. Additionally, Kim et al. (2023) and Stelmakh et al. 160 (2022) aim to generate answers that encompass all possible responses. In this paper, we focus on 161 generating clarification questions using both pre-trained and fine-tuned language models. Further-

Ambiguity Type	Example
Lexical (19.65%)	Q: Why did the medieval church need an alarm arrangement?
	C: Are you asking if the medieval church needed an alarm arrangement to wake people up or to signal a threat? CA: To wake people up / QA: The medieval church used an alarm arrangement to wake people up. CA: To signal a threat / QA: The medieval church used an alarm arrangement to signal a threat to the community.
Syntactic (15.67%)	Q: What aspect of creating new roles would most weaken the limited impact thesis criticized by women's rights activists?
	C. Do you mean new roles in high-tech and service sectors or the roader societal and economic transitions? C.A: New roles in high-tech and service. (Q.A: The aspect where new roles in the high-technology and service sectors were being created. C.A: Societal and economic transitions. / Q.A: Critics argue the transition was painful but temporary for broader societal impact.
Incomplete (30.46%)	Q: What is found inside living organisms according to the initial segments?
	C. Are you asking about the main chemical elements present in hving organisms or the organic molecules found whimit them? C.A: The main chemical elements present in living organisms. / Q.A: Oxygen, carbon, hydrogen, and nitrogen make up most of living organisms' mass. C.A: The organic molecules found within living organisms. / Q.A: Oxygen, carbon, hydrogen, and nitrogen make up most of living organisms' mass.
Contextual (34.22%)	C. What insishis were sained about their twing organisms. (Q-A) organic molecules are essentiar to twing organisms. O: What insishis were sained about their twing organisms. (Q-A) organic molecules are essentiar to twing organisms.
Contextual (34.22.10)	C. Are you asking about the communication function or the seasonality and traits of forehead rubbing? C.A: The communication function / Q.A: Studies show forehead rubbing in deer communicates identity, sex, and dominance through scent. C.A: The seasonality and howsical traits / O.A: Studies show forehead rubbing in deer communicates identity, sex, and dominance through scent.

Table 3: Four types of ambiguity: Lexical, Syntactic, Incomplete, and Contextual Ambiguity. In the examples, 'Q' denotes an ambiguous question, 'C' represents a clarification question, 'C_A' stands for the clarification answer, and 'Q_A' signifies the answer to the ambiguous question after clarification.

more, we evaluate the performance of large language models (LLMs) in the context of the answers
 to the generated clarification questions.

179 **Datasets on Ambiguity** Many datasets address ambiguity in conversation and question answering. 180 To our knowledge, Braslavski et al. (2017) introduces the first ambiguous dataset using community 181 question-answering websites. Rao & Daumé III (2018) utilize data from StackExchange, while 182 Saeidi et al. (2018) focus on rules and laws. Wu et al. (2023) creates an ambiguous dataset by 183 extracting conversations from Wikipedia using web searches. Other ambiguous datasets are based 184 on well-known public datasets. For example, Guo et al. (2021) propose Abg-CoQA, which clarifies 185 ambiguities based on CoQA (Reddy et al., 2019). Min et al. (2020) generate AmbigQA for opendomain question answering based on NQ-Open (Lee et al., 2019). Stelmakh et al. (2022) uses AmbigQA to enhance long-form QA in the context of ambiguity. Lee et al. (2023) further refine 187 AmbigQA with the assistance of InstructGPT (Ouyang et al., 2022). 188

Though there are many ambiguous datasets, these datasets do not provide annotations of different types of ambiguity. Besides, most datasets consider using data from simple areas like community conversations or public datasets which may used to train the language models. On the other hand, none of the previous works consider evaluating the close-source commercial LLMs such as GPTo1. To this end, Our dataset contains not only annotations of ambiguous types but also high-quality passages and questions from various science areas. We include the evaluation of our datasets on the commercial LLMs as well, which distinguishes our work from previous.

196 197

198 199

200

201

202

3 DATASET COLLECTION

We build Abg-SciQA based on various questions and different areas in both natural and social science. We collect stories, questions, and corresponding answers from TOEFL, IELTS, GRE, and GMAT reading comprehensive Mock Tests. We will start introduction from tasks in Abg-SciQA.

- 203204 3.1 TASK DEFINITION
- 205

Given a story S and a question-answer pair $\{Q, A\}$, the task is to resolve any ambiguity. We consider four unique tasks and three of them have been considered in previous work Guo et al. (2021) while ambiguity type classification is the new task we introduce in our paper.

Ambiguity Detection: Determining whether the question Q is ambiguous based on the given story S and the question-answer pair $\{Q, A\}$.

Ambiguity Type Classification: Given four definitions and examples of ambiguity, LLMs are asked to predict the type of ambiguity of the ambiguous question Q.

Clarification Question Generation: If the question Q is found to be ambiguous, generate a clarification question CQ that will help future disambiguate Q and make the answer of ambiguous question more specific.



Figure 3: The Abg-SciQA pipeline for generating ambiguous questions. It starts with an exam story and question, along with a predefined ambiguity type, to generate an ambiguous question and two answers using GPT-40. The story is then revised to align with the ambiguous answers, followed by generating clarification questions. Finally, Claude-3.5-Sonnet evaluates ambiguity and answer consistency. Only valid entries are added to the dataset, with a subset undergoing human evaluation before finalizing Abg-SciQA.

Clarification-Based Question Answering: Using the story S, the original question-answer pair $\{Q, A\}$, the generated clarification question CQ, and a possible response R to CQ, provide a clear and unambiguous answer to the original question Q. Different responses R should lead to different answers A_R .

245 246

3.2 MATERIAL COLLECTION PROCESS

247 The previous datasets on ambiguity are mainly based on the public Natural Language Processing 248 (NLP) dataset. For example, Abg-CoQA (Guo et al., 2021) is generated based on CoQA (Reddy 249 et al., 2019) and AmbigQA (Min et al., 2020), generated based on NQ-open (Lee et al., 2019). Most 250 of these public datasets are based on some simple tests and CoQA is based on the children's stories 251 from MCTest (Richardson et al., 2013) and middle and high school English exams from RACE (Lai et al., 2017). These exams are less challenging compared with more advanced science questions 253 and thus may be easier for LLMs to understand the contexts. Therefore, in this paper, we tend to 254 use passages and questions from advanced exams like TOEFL Mock Test. To collect the necessary 255 material, we gather stories (passages), questions, and answers in various nature and social science areas from the Internet. However, most of these questions are not ambiguous. Therefore, we need 256 to generate ambiguous questions. 257

258 259

3.3 Ambiguous Question Generation

260 We provide the generation and evaluation pipeline in Fig. 3. Since most questions in these tests 261 are clear and not ambiguous, we employ GPT-40 (Achiam et al., 2023) to automatically generate 262 ambiguous questions and corresponding answers based on unambiguous questions, following estab-263 lished definitions of ambiguity. To facilitate this process, we design a prompt for the LLMs in Fig. 5 264 that includes definitions of four types of ambiguity, each accompanied by examples. This approach 265 helps LLMs understand the precise meaning of ambiguity, as human perceptions may differ from 266 those of the models. By providing clear definitions, we reduce the risk of misunderstandings. For instance, different samples in Abg-CoQA (Guo et al., 2021) may be interpreted inconsistently re-267 garding continuity, as shown in Fig. 2. Two similar cases may present different types of ambiguity 268 due to this issue, potentially leading to unfair evaluations. We only prompt GPT-40 to generate one 269 type of ambiguity one time.

270 What's more, to ensure that the ambiguous answers are supported by the story. We continue using 271 GPT-40 to revise the original story so that it aligns with the generated answers. The detailed prompt 272 for story revision is in Fig. 6. In addition to producing ambiguous questions, we also ask GPT-40 273 to provide clarification questions aimed at resolving the ambiguity, and the prompt for clarification 274 generations is in Fig. 7.

- 276 3.4 AMBIGUITY TYPES
- 277

275

278 Table 2 shows a formal definition of ambiguity types in our paper and Table 3 shows a breakdown and examples of the ambiguity type in Abg-SciOA. We define ambiguity with four categories : 279 Lexical Ambiguity, Syntactic Ambiguity, Incomplete Ambiguity, and Contextual Ambiguity. 280

281 Lexical Ambiguity occurs when a word has multiple meanings or multiple interpretations. For ex-282 ample, in Table 3, the Lexical Ambiguity example is: "Why did the medieval church need an alarm 283 arrangement?". "The alarm arrangement" can be interpreted as waking people up before a scheduled event or a signal to signal a threat or danger to the community. Therefore, the answer will vary 284 depending on the context. 285

286 Syntactic Ambiguity arises from the structure of a sentence, allowing multiple interpretations de-287 pending on how it is parsed. For example, in Table 3, the Syntactic Ambiguity example is: "What 288 aspect of creating new roles would most weaken the limited impact thesis criticized by women's 289 rights activists?" The ambiguity stems from the phrase "criticized by women's rights activists," which can modify either "limited impact thesis" or "creating new roles." As a result, the sentence 290 can yield two distinct interpretations depending on its structural parsing. 291

292 Incomplete Ambiguity occurs when a statement or question lacks essential contextual information, 293 such as location, time, event, or people, resulting in multiple possible interpretations. This type 294 of ambiguity arises from the practical understanding of language in context. For example, in Table 3, the question demonstrating Incomplete Ambiguity is: "What is found inside living organisms 295 according to the initial segments?" In this case, the question lacks specificity, making "What" am-296 biguous—it could refer to the main chemical elements or organic molecules. Consequently, two 297 possible answers arise based on different interpretations due to the absence of essential information. 298

299 Contextual Ambiguity occurs when a question is clear and unambiguous but contains two possible 300 answers due to differing contexts, interpretations, or sources. For example, in Table. 3, the Contex-301 tual Ambiguity question is: "What insights were gathered from the research on how deer use their foreheads when they rub against trees?" There are two possibles insights based on the context: one 302 is the communication function of forehead rubbing. Another is the seasonality and physical traits 303 associated with the behavior. 304

305 Each of these types of ambiguity represents distinct challenges in understanding and interpreting 306 ambiguity. These four categories cover a broad range of ambiguous cases. 307

- 308
- 309

3.5 QUALITY CONTROL

To ensure the quality of our results, we use a two-stage evaluation. In the first stage, we use Claude-310 3.5-Sonnet to assess the quality of the generated dataset. Specifically, we provide Claude-3.5-Sonnet 311 with the ambiguous dataset, a corresponding ambiguity type definition, and the adjusted story. We 312 then ask whether the proposed question meets the ambiguity requirements and aligns with the given 313 ambiguity definition. Additionally, we ask Claude-3.5-Sonnet whether the answers to the ambigu-314 ous and clarification questions can be found in the revised story to ensure their accuracy. In total, we 315 generate 25,000 questions using GPT-40, and Claude filters out 12,502 entries. In the second stage, 316 we conduct human evaluation to verify that Claude's assessment aligns with human judgment. The 317 human evaluators are native English speakers from various universities in the U.S. We begin by con-318 ducting a human evaluation to identify and remove poorly constructed ambiguous questions from the 319 dataset and human filtered out around 2,296 entries. After this refinement, We also randomly select 320 50 entries from the dataset and have four human evaluators answer the same questions as Claude 321 to assess its reliability. To ensure the quality of human evaluation, evaluators must pass an exam consisting of five entries labeled by a human expert before proceeding with the assessment. We use 322 Cohen's Kappa statistic (ML. et al., 2012) to measure the agreement between Claude's evaluations 323 and human evaluations, yielding a final result of 0.6535, which indicates moderate agreement between Claude's assessments and human evaluations. This suggests that Claude's assessments are generally reliable.
 326

3.6 DATA ANALYSIS

Abg-SciQA is composed of 1,353 stories and 13,729 330 questions, where 10,202 questions are annotated as ambiguous. The comparison in Table 1 shows that 331 our dataset if one of the largest datasets in the am-332 biguity area. We present the distributions of source 333 domains for Abg-SciQA in Fig. 4. Abg-SciQA con-334 tains more than 10 different areas from either natural 335 science or social science questions. In Table 3, we 336 provide the distribution of different ambiguity types 337 in Abg-SciQA. Our diverse types of areas and types 338 of ambiguity can lead to a better evaluation and im-339 provement of LLMs' ability to solve ambiguity.



Figure 4: The distribution of story domains in Abg-SciQA. Among all domains, History questions account for the most in Abg-SciQA, followed by Biology.

341 342 4 EVALUATION ON DATASET

In this paper, we evaluate Abg-SciQA in both closed-source commercial LLMs and open-source
public LLMs to show how different language models deal with our dataset comprehensively. We
also include the results of fine-tuned LLMs with Abg-SciQA to show Abg-SciQA can guide the
improvement of LLMs. All of our experiments are done on one single NVIDIA A100-80G GPU. In
detail, we randomly sample 80% of Abg-SciQA as the training set and use the rest as the evaluation
set. We prompt the LLM using few-shot.

349 350

351

340

327

328

4.1 EVALUATION METRICS

Abg-SciQA contains four tasks, one more than previous work (Guo et al., 2021). For Ambiguity 352 Detection, since we treat this task as a binary classification, we report precision, recall, and F1 as 353 the evaluation metrics. For Ambiguity Type Classification, since this task can be treated as a multi-354 class classification, we compute macro-precision, macro-recall, and macro-F1. For Clarification 355 Question Generation, we use BLEU and Rouge-L as metrics with the labeled clarification question 356 as the gold standard. In addition to directly measuring the quality through automatic metrics, we 357 also manually evaluate whether the generated question is reasonable and helpful for clarifying the 358 existing ambiguity for a small subset. Finally, for Clarification-Based Question Answering, we 359 follow the common practice to compute the macro-average F1 score of word overlap (Reddy et al., 360 2019).

361 362

363

4.2 EVALUATION ON CLOSED-SOURCE LLMS

We evaluate Abg-SciQA in 3 closed-source LLMs: 1) GPT-01 (Achiam et al., 2023) 2) Gemini (Team et al., 2023) 3) Claude (Anthropic, 2024). We do not evaluate Abg-SciQA on GPT-40 (Achiam et al., 2023) because we use GPT-40 (Achiam et al., 2023) to generate our dataset. We provide the results in Table 4 and detailed results for each type of ambiguity in the Table 5. We have the following observations given the results:

1)Different closed-source LLMs exhibit varying performance across different tasks. For example, GPT-01 performs the best in Ambiguity Detection, while Claude-sonnet shows a stronger performance in the task like Ambiguity Type Classification.

2) However, as we can see, even Claude-sonnet and GPT-o1 do not solve the problem very well.
For example, the F1 score for Ambiguity Detection in Claude-sonnet is 0.5077, and the F1 score for
Ambiguity Type Classification in GPT-o1 is only 0.2194. What's more, all the model performances
is really bad on Clarification-Based QA. This indicates that our dataset is highly challenging, and
even state-of-the-art models struggle to handle it well.

376 3) Compared to all tasks, Ambiguity Detection appears to be the easiest. This is likely because we
 377 provide a very clear definition of each type of ambiguity along with corresponding examples, making it easier for LLMs to detect ambiguity in a question. However, if we rely solely on the model to

³⁷⁸ classify the type of ambiguity without these aids, the performance drops significantly.

4) Even though Ambiguity Type Classification is a comparatively easier problem among all tasks
provided in Abg-SciQA, the results in Table 5 show that it is hard for LLMs to understand all types
of ambiguity. Claude-sonnet with the highest overall performance on Ambiguity Type Classification shows a very good understanding of Lexical Ambiguity and Contextual Ambiguity. However,
Claude-sonnet can hardly understand the rest of two types.

	Ambiguity Detection Type Classification Clarification Generation				ion Generation	QA			
Model	Precision	Recall	F1	Precision	Recall	F1	BLEU	Rouge-L	F1
GPT-o1-mini	0.6707	0.1375	0.2282	0.2039	0.2425	0.1764	0.1393	0.3517	0.0062
GPT-01	0.7985	0.5450	0.6478	0.2548	0.2550	0.2194	0.1490	0.3516	0.0025
Gemini-1.0	0.9259	0.0312	0.1605	0.3292	0.3325	0.3013	0.0716	0.2433	0.0012
Gemini-1.5	0.9685	0.2312	0.3733	0.4004	0.3112	0.2106	0.0822	0.2585	0.0050
Claude-haiku	0.8703	0.1762	0.2931	0.6263	0.3737	0.3416	0.0752	0.2591	0.0025
Claude-sonnet	0.3675	0.8212	0.5077	0.6818	0.4587	0.3992	0.1054	0.2933	0.0037

Table 4: Performance of different closed-source LLMs for all tasks provided by Abg-SciQA. We highlight the **best** performance and the <u>second best</u>. The results show that Claude-sonnet performs the best. However, even a powerful model like Claude-3.5-sonnet and GPT-o1 still perform not very well in all tasks. The full model name is o1-mini-2024-09-12, o1-preview-2024-09-12, gemini-1.0-pro, gemini-1.5-pro, claude-3-5-haiku-20241022, and claude-3-5-sonnet-20241022.

Model	Lexical	Syntactic	Incomplete	Contextual
GPT-o1-mini	0.2461	0.0417	0.0515	0.3666
GPT-o1	0.2152	0.1107	0.1905	0.3613
Gemini-1.0	0.5298	0.0000	0.3521	0.3236
Gemini-1.5	0.3740	0.0000	0.4222	0.0465
Claude-haiku	0.7003	0.0392	0.2121	0.4150
Claude-sonnet	0.8177	0.0583	0.2472	0.4737

Table 5: Detailed F1 score of Ambiguity Type Classification for each type on closed-source LLMs.
 The results show that LLMs often cannot understand ambiguity well even though the state-of-art like GPT-o1 and Claude-sonnet.

409 4.3 EVALUATION ON OPEN-SOURCE LLM

384

392

393

394

395

Now We evaluate Abg-SciQA in 5 open-source LLM: 1) Llama2 (Touvron et al., 2023) 2) Llama3 3)
Gemma (Team et al., 2024) 4) Phi3 (Abdin et al., 2024) 5)Mistral (Albert et al., 2023). We present our overall results in Table 6 and Table 7. We have the following observations:

1) In Table 6, similar to previous results with closed-source LLMs, all open-source models struggle to solve ambiguous problems effectively. Additionally, compared to closed-source LLMs, some
open-source models perform worse, likely due to the significant difference in the number of parameters between closed-source and open-source models.

2) In Table 6, even though Mistral-0.1 achieves the best performance in Ambiguity Detection and Ambiguity Type Classification, it can only understand three types of ambiguity well, as shown in Table 7.

3)Aside from Mistral-0.1 and Gemma-2B, the other open-source LLMs perform poorly across most tasks. While some models, such as Llama3-1.8B, excel in Ambiguity Detection, their performance in Ambiguity Type Classification, Clarification Generation, and Clarification-Based Question Answering remains subpar. The variation in performance across tasks suggests that these models struggle to consistently address ambiguity, highlighting the need for further improvement and fine-tuning.

Finally, we tend to fine-tune the open-source LLMs with Abg-SciQA to see if our dataset can guide
the further training of LLMs. In detail, we randomly sample 80% of Abg-SciQA as the training set
and use the rest as the evaluation set. We choose the same open-source models in the previous section
as training models. We use LoRA to train all the models with AdamW (Loshchilov & Hutter, 2017),
Lora rank 8, learning rate 5e-6, and training epochs 3. And we present the results of the evaluation
set in Table 8 and Table 9. We can have the following observations:

1) Fine-tuning can significantly improve the performance of chosen open-sourced LLMs. We can see that the best model Llama2-13B can beat the closed-source LLMs in Ambiguity Detection and

	Amb	Ambiguity Detection			Type Classification			Clarification Generation		
Model	Precisior	n Recall	F1	Precision	Recall	F1	BLEU	Rouge-L	F1	
Llama2-7	B 0.7528	0.2219	0.2731	0.2958	0.1200	0.1655	0.0303	0.1714	0.0096	
Llama2-1	3B 0.7491	0.3421	0.2527	0.2430	0.1737	0.1666	0.0050	0.0912	0.0083	
Llama3.1	-8B 0.8837	0.3612	0.5128	0.3453	0.0662	0.0969	0.0035	0.0741	0.0175	
Llama3.2	-3B 0.8850	0.0962	0.1736	0.2460	0.2062	0.1338	0.0027	0.0658	0.0187	
Gemma-2	2B 0.7689	0.5825	0.6628	0.2839	0.2714	0.1668	0.0069	0.0821	0.0054	
Phi3.5	0.8089	0.5187	0.6321	0.3071	0.2525	0.1415	0.0326	0.1627	0.0154	
Mistral-0	.1 0.7940	0.7348	0.7632	0.2615	0.2450	0.2240	0.0482	0.1984	0.0112	
Mistral-0	.2 0.8235	0.1925	0.3120	0.3821	0.2587	0.1564	0.0697	0.2300	0.0087	

Table 6: Performance of different open-source LLMs for all tasks provided by Abg-SciQA. we highlight the **best** performance and the <u>second best</u>. The results show that open-source LLMs are good at Ambiguity Detection. However, most of them fall short on other tasks, compared with more powerful closed-source LLMs.

Model	Lexical	Syntactic	Incomplete	Contextual
Llama2-7B	0.2202	0.0667	0.1479	0.2276
Llama2-13B	0.2395	0.0000	0.3983	0.0287
Llama3.1-8B	0.0857	0.0000	0.1233	0.1787
Llama3.2-3B	0.3611	0.0747	0.0538	0.0457
Gemma-2B	0.1812	0.0291	0.0095	0.1573
Phi3.5	0.0294	0.0000	0.1379	0.3987
Mistral-0.1	0.2717	0.0611	0.3088	0.2544
Mistral-0.2	0.0100	0.0000	0.2242	0.3916

Table 7: Detailed F1 score of Ambiguity Type Classification for each type on open-source LLMs without fine-tuning. The results show that most LLMs cannot understand ambiguity well before fine-tuning.

460 Ambiguity Type Classification, demonstrating the effectiveness of training on Abg-SciQA to solve the problem.

461 the problem.
2) Fine-tuning significantly improved performance for Ambiguity Detection. However, similar to open-source LLMs, there still remain a lot of problems for LLMs in understanding the meaning of all types of ambiguity. For example, Though Llama2-13B has a good performance on Ambiguity Type Classification after fine-tuning, Llama2-13B can only have a good understanding of two types of ambiguity in Table 9. These results indicate that further improvement is needed and Abg-SciQA can help to guide development.

	Ambiguity Detection			Туре	Type Classification			Clarification Generation	
Model	Precision	Recall	F1	Precision	Recall	F1	BLEU	Rouge-L	F1
Llama2-7B	0.9989	0.9989	0.9989	0.4741	0.5055	0.6000	0.1847	0.4773	0.0114
Llama2-13B	0.9982	0.9948	0.9969	0.7812	0.7894	0.7917	0.1224	0.4256	0.0167
Llama3.1-8B	0.9683	0.9739	0.9829	0.6650	0.7285	0.5235	0.1749	0.2316	0.0195
Llama3.2-3B	1.0000	0.9692	0.9843	0.5459	0.6977	0.6451	0.0920	0.3736	0.0187
Gemma-2B	1.0000	0.9794	0.9896	0.4990	0.4269	0.6287	0.0313	0.1928	0.0092
Phi3.5	0.9979	0.9984	0.9984	0.5196	0.5502	0.5948	0.0738	0.2018	0.0179
Mistral-0.1	0.9858	0.9986	0.9923	0.5107	0.6155	0.6235	0.0759	0.2935	0.0141
Mistral-0.2	0.9979	0.9983	<u>0.9984</u>	0.5069	0.5218	0.6297	0.0775	0.2821	0.0093

Table 8: Performance of different open-source LLMs for all tasks after fine-tuning on Abg-SciQA. we highlight the **best** performance and the <u>second best</u>. The results show that fine-tuning on Abg-SciQA can significantly increase the performance of all LLMs and make smaller-size models become even better than closed-source LLMs.

482 4.4 TRANSFER ABILITY FOR ABG-SCIQA

We evaluated open-source LLMs that are fine-tuned on Abg-SciQA by Abg-CoQA (Guo et al., 2021) to see whether Abg-SciQA can help to increase the ability to solve other general ambiguity questions. We focused on Abg-CoQA because Abg-CoQA has more comprehensive tasks (three

36	Model	Lexical	Syntactic	Incomplete	Contextual
7				· · · ·	
0	Llama2-7B	0.8704	0.0098	0.0000	0.6881
)	Llama2-13B	0.9684	0.7232	0.5350	0.8121
	Llama3.1-8B	0.9386	0.6792	0.0613	0.7529
	Llama3.2-3B	0.9635	0.0744	0.1921	0.7079
	Gemma-2B	0.9804	0.0000	0.0000	0.6982
	Phi3.5	0.9150	0.0952	0.1373	0.6801
	Mistral-0 1	0 9798	0.0288	0.0479	0 6939
	Mistral 0.2	0.0025	0.0200	0.0000	0.6071
	Wilstrai-0.2	0.9925	0.0280	0.0000	0.09/1

Table 9: Detailed F1 score of Ambiguity Type Classification for each type on open-source LLMs after fine-tuning. The results show that most LLMs can only understand two type well even after fine-tuning.

tasks in total). To better analyze transfer ability, we consider the performance increasing between models without fine-tuning and models with fine-tuning on Abg-SciQA. Our results are presented in Table 10. Based on these results, we have the following observations:

1) Compared to evaluating open-source LLMs, Abg-CoQA performed better on our fine-tuned model across tasks. For instance, in Ambiguity Detection with Llama2-7b, the F1 score has an improvement of approximately two times compared to the model without fine-tuning.

2) When evaluating on Abg-SciQA, Llama3.2-3B shows minimal performance improvement after
fine-tuning. However, when evaluated on Abg-CoQA, its performance improves significantly in recall and F1, particularly in Ambiguity Detection. While its overall improvement is not the highest
among all models, these results further validate that fine-tuning on Abg-SciQA enhances models'
ability to handle a wide range of ambiguous questions, demonstrating its effectiveness as a training
set for ambiguity resolution.

	Ambiguity Detection			Clarificatio	QA	
Model	Δ Precision	Δ Recall	Δ F1	Δ BLEU	Δ Rouge-L	Δ F1
Llama2-7B	-2.7%	388.2%	100.2%	865.3%	432%	1.8%
Llama2-13B	53.9%	-35.3%	-6%	276.9%	615.1%	30.3%
Llama3.1-8B	-15.1%	358.4%	14%	100.5%	245.9%	1.2%
Llama3.2-3B	-1%	586.7%	114.9%	53.9%	157.7%	7.3%
Gemma-2B	2.1%	5.1%	1.8%	27.9%	40.3%	31.5%
Phi3.5	4.9%	-0.8%	10%	60.1%	33.8%	30.4%
Mistral-0.1	2.5%	30.5%	-0.6%	205.5%	214.6%	15.2%
Mistral-0.2	-2.7%	-12.6%	-24.1%	116.7%	174.7%	21.2%

Table 10: Performance increasing of different open-source LLMs for all tasks on Abg-CoQA after fine-tuning on Abg-SciQA. we highlight the **best** performance and the <u>second best</u>. The results show that in general, fine-tuning on Abg-SciQA can help to improve the performance on Abg-CoQA, indicating a good generalization ability of Abg-SciQA.

5 CONCLUSION

In this paper, we introduce Abg-SciQA, a dataset aiming at evaluating LLMs on detecting and solv-ing ambiguity comprehensively. Derived from advanced science questions and enhanced with gen-erated ambiguous questions, Abg-SciQA encompasses four key tasks to analyze ambiguity better. Our extensive experiments on both closed-source and open-source LLMs reveal that even state-of-the-art models struggle with these tasks, highlighting areas for improvement. Notably, fine-tuning open-source LLMs on Abg-SciQA leads to substantial performance gains, demonstrating its po-tential to guide LLM development in ambiguity handling. Additionally, we evaluated Abg-CoQA using Abg-SciQA fine-tuned models, which also showed significant improvement. This demonstrates the flexibility of Abg-SciQA fine-tuned models and its potential to perform well on other datasets. Abg-SciQA thus serves as a valuable benchmark for advancing language understanding on ambiguous questions.

540 REFERENCES

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany
 Awadalla, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023, 2023.
- Jiang Albert et al. Mistral 7b. *arXiv preprint arXiv:2310.06825.*, 2023.
- Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and W Bruce Croft. Asking clarifying questions in open-domain information-seeking conversations. In *Proceedings of the 42nd international acm sigir conference on research and development in information retrieval*, pp. 475–484, 2019.
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeffrey Dalton, and Mikhail Burtsev.
 Building and evaluating open-domain dialogue corpora with clarifying questions. *arXiv preprint arXiv:2109.05794*, 2021.
- A.I. Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- Pavel Braslavski, Denis Savenkov, Eugene Agichtein, and Alina Dubatovka. What do you mean
 exactly? analyzing clarification questions in cqa. In *Proceedings of the 2017 conference on conference human information interaction and retrieval*, pp. 345–348, 2017.
- Herbert H Clark and Susan E Brennan. Grounding in communication. American Psychological Association, 1991.
- Yang Deng, Wenqiang Lei, Lizi Liao, and Tat-Seng Chua. Prompting and evaluating large lan guage models for proactive dialogues: Clarification, target-guided, and non-collaboration. *arXiv preprint arXiv:2305.13626*, 2023a.
- Yang Deng, Shuaiyi Li, and Wai Lam. Learning to ask clarification questions with spatial reasoning. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2113–2117, 2023b.
- Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. Evaluating large language models in class-level code generation. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pp. 1–13, 2024.
- Zachary Eberhart and Collin McMillan. Generating clarifying questions for query refinement in source code search. In 2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER), pp. 140–151. IEEE, 2022.
- Meiqi Guo, Mingda Zhang, Siva Reddy, and Malihe Alikhani. Abg-coqa: Clarifying ambiguity in conversational question answering. *In 3rd Conference on Automated Knowledge Base Construction*, 2021.
- Baber Khalid, Malihe Alikhani, and Matthew Stone. Combining cognitive modeling and reinforce ment learning for clarification in dialogue. *In Proceedings of the 28th International Conference on Computational Linguistics*, 2020.
- Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joonsuk Park, and Jaewoo Kang. Tree of clarifications: Answering ambiguous questions with retrieval-augmented large language models. *arXiv* preprint arXiv:2310.14696, 2023.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale ReAding
 comprehension dataset from examinations. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel
 (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Process- ing*, pp. 785–794, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082. URL https://aclanthology.org/D17-1082.

- Dongryeol Lee, Segwang Kim, Minwoo Lee, Hwanhee Lee, Joonsuk Park, Sang-Woo Lee, and Kyomin Jung. Asking clarification questions to handle ambiguity in open-domain qa. *arXiv preprint arXiv:2305.13808*, 2023.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6086–6096, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1612. URL https://www.aclweb.org/anthology/P19-1612.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint
 arXiv:1711.05101, 2017.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. Ambigqa: Answering ambiguous open-domain questions. *arXiv preprint arXiv:2004.10645*, 2020.
- McHugh ML. et al. Interrater reliability: the kappa statistic. *Biochem Med (Zagreb)*, 2012.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Matthew Purver, Jonathan Ginzburg, and Patrick Healey. On the means for clarification in dialogue.
 Current and new directions in discourse and dialogue, pp. 235–255, 2003.
- Sudha Rao and Hal Daumé III. Learning to ask good questions: Ranking clarification questions
 using neural expected value of perfect information. *arXiv preprint arXiv:1805.04655*, 2018.

- Siva Reddy, Danqi Chen, and Christopher D. Manning. Coqa: A conversational question answering
 challenge. *Transactions of the Association for Computational Linguistics*, 2019.
- Matthew Richardson, Christopher J.C. Burges, and Erin Renshaw. MCTest: A challenge dataset for
 the open-domain machine comprehension of text. In David Yarowsky, Timothy Baldwin, Anna
 Korhonen, Karen Livescu, and Steven Bethard (eds.), *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 193–203, Seattle, Washington, USA,
 October 2013. Association for Computational Linguistics. URL https://aclanthology.
 org/D13-1020.
- Kepa Joseba Rodríguez and David Schlangen. Form, intonation and function of clarification requests
 in german task-oriented spoken dialogues. *the 8th workshop on the semantics and pragmatics of dialogue*, 2004.
- Marzieh Saeidi, Max Bartolo, Patrick Lewis, Sameer Singh, Tim Rocktäschel, Mike Sheldon, Guillaume Bouchard, and Sebastian Riedel. Interpretation of natural language rules in conversational machine reading. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2087–2097, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1233. URL https://aclanthology.org/D18-1233.
- David Schlangen. Causes and strategies for requesting clarification in dialogue. *the 5th SIGdial Workshop on Discourse and Dialogue at HLT-NAACL 2004*, pp. 136–143, 2004.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. Asqa: Factoid questions meet long-form answers. *arXiv preprint arXiv:2204.06092*, 2022.
- 642 Svetlana Stoyanchev, Alex Liu, and Julia Hirschberg. Towards natural clarification questions in
 643 dialogue systems. *In AISB symposium on questions, discourse and dialogue*, 20, 2014.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, et al. Gemini: a family of highly capable multi modal models. *arXiv preprint arXiv:2312.11805*, 2023, 2023.
- 647 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295, 2024, 2024.*

648 649 650	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288, 2023</i> , 2023.
651 652 653 654	Zeqiu Wu, Ryu Parish, Hao Cheng, Sewon Min, Prithviraj Ammanabrolu, Mari Ostendorf, and Han- naneh Hajishirzi. Inscit: Information-seeking conversations with mixed-initiative interactions. <i>Transactions of the Association for Computational Linguistics</i> , 11:453–468, 2023.
655 656 657	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. <i>arXiv preprint arXiv:2306.13063</i> , 2023.
658 659 660	Samuel Yu, Shihong Liu, Zhiqiu Lin, Deepak Pathak, and Deva Ramanan. Language models as black-box optimizers for vision-language models. <i>arXiv preprint arXiv:2309.05950</i> , 2023.
661 662 663	Hamed Zamani, Susan Dumais, Nick Craswell, Paul Bennett, and Gord Lueck. Generating clarify- ing questions for information retrieval. In <i>Proceedings of the web conference 2020</i> , pp. 418–428, 2020.
664 665 666	Michael JQ Zhang and Eunsol Choi. Clarify when necessary: Resolving ambiguity through interaction with lms. <i>arXiv preprint arXiv:2311.09469</i> , 2023.
667	
668	
669	
670	
670	
672	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
090 607	
605	
600	
700	
701	

702 A DATASET AND CODE

As mentioned in the abstract, our code and dataset can be found at https://anonymous. 4open.science/r/Abg-Sci-DF10/README.md.

B PROMPT FOR GENERATING AMBIGUOUS QUESTION

We present our prompt for generating ambiguous questions in Fig. 5. This prompt includes the story, original question, and original answer. It then defines and provides examples of four types of ambiguity before instructing GPT-40 to generate the required output using chain-of-thought reasoning. Any samples that do not conform to the expected format are adjusted by human annotators.

714 In Fig. 6, we show the prompt used for story revision, ensuring that the ambiguous question aligns 715 more closely with the story. Fig. 7 details the prompt for generating clarification questions. Fi-716 nally, Fig. 8 presents the prompts given to Claude-sonnet-3.5 to assess the alignment of ambiguous 717 questions with the ambiguity definitions and verify the correctness of all answers based on the story.

718

704

705

706

708

719

720	Ambiguous Generation
721	Story: {story}, Original Question: {original question}, Original Answer: {original answer}.
722	The ambiguous question has four types: Lexical, Syntactic, Incomplete, and Contextual Ambiguity
703	Example1: Lexical Ambiguity Question: {definition and a question example}
723	clarification question: {clarification question},
724	Example 2: Syntactic Ambiguity Question: { definition and a question example }
725	clarification question: {clarification question},
726	Example3: Incomplete Ambiguity Question: { definition and a question example }
727	clarification question: {clarification question}.
728	Example4: Contextual Ambiguity Question: { definition and a question example }
700	clarification question: {clarification question}.
729	Instructions:
730	Please use the story, original question, and original answer to generate an ambiguous question based on <lexical,< th=""></lexical,<>
731	Syntactic, Incomplete, or Contextual> Ambiguity. Please think step by step and tell me the reason why the
732	question you generated is ambiguous and give me two possible answers based on the story. Please generate the
733	ambiguous question based on the story, original question, and original answer rather than the examples. The answers of ambiguous questions must be clearly found in the story and please give me two ambiguous answers
734	Output Format:
735	Ambiguous Question: <your generated="" question=""></your>
736	Ambiguous Answer 1: <first answer="" possible=""></first>
737	Ambiguous Answer 2: <second answer="" possible=""></second>
700	Explanation: <explanation ambiguous="" is="" of="" question="" the="" why=""></explanation>
130	
739	

Figure 5: The prompt for ambiguous question generations in Abg-SciQA.

744 745

746

C DATASET STRUCTURE

In this section, we provide the structure of Abg-SciQA. In detail, Abg-SciQA is stored in a JSON file, and in Fig. 9 shows the detailed structure of the JSON file.

747 748 749

D COMPARISON OF ABG-SCIQA ON AMBIGUITY DETECTION

In this section, we present both a successful and an unsuccessful case of Ambiguity Detection, comparing the performance of a model before and after fine-tuning. As shown in Fig. 10, we use
Llama3.1-8B to perform Ambiguity Detection, both in its original state and after fine-tuning with Abg-SciQA. The results demonstrate that the fine-tuned model effectively identifies ambiguity. However, in Fig. 11, even after fine-tuning, the model fails to detect ambiguity, highlighting its limitations in certain cases.

756	
757	Story Revision
758	Story: {story}, Ambiguous Question: {generated question}, Ambiguous Answer 1: { answer for ambiguous
759	is ambiguous answer 2: {answer for ambiguous question}, Explanation: {reasons why the question
760	Instructions:
761	Please revise the story based on the ambiguous question, ambiguous answers, and explanation, and make the
762	ambiguous answer true. Please give me the full story after revised. You should make sure the ambiguous answer is
763	followed the revised story. The ambiguous question and answers must be the same as the input. The Ambiguous answer cannot be Partially Correct. It should be fully correct based on the revised story.
764	Output Format:
765	Revised Story: <revised story=""></revised>
766	Ambiguous Question: <ambiguous question=""></ambiguous>
707	Ambiguous answer1: <answer ambiguous="" for="" question=""> Ambiguous answer2:<answer ambiguous="" for="" question=""></answer></answer>
760	Antorguous unswerz, canswer for antorguous questions
709	
771	Figure 6: The prompt for story revision in Abg-SciQA.
772	
773	Clarification Generation
774	Revised Story (revised story) Ambiguous Question (ambiguous question) Ambiguous answer! (answar for
775	ambiguous question}, Ambiguous answer2:{answer for ambiguous question}
776	Instructions:
777	Please ask a clarification question to clarify the ambiguous question based on the revised story, ambiguous
778	question, and answers. For the c_answer, please don't start with 'If you are referring to'. If the Clarification
779	should be I'm asking about B. The Ambiguous Answers should not be the same as Clarification Answers Please
780	provide me with the clarification answer in format: "
781	Output Format:
782	Clarification Question: <clarification question=""> c_answer1:<clarification answer1=""></clarification></clarification>
783	c_answer2: <clarification answer2=""></clarification>
784	
785	Figure 7: The prompt for clarification question generation in Abg-SciOA
786	rigure 7. The prompt for charmention question generation in 1.65 berg.r.
787 788 789	E COMPARISON OF ABG-SCIQA ON AMBIGUITY TYPE CLASSIFICATION
790 791 792 793	in this section, we examine both a successful and an unsuccessful case of Ambiguity Type Classification, comparing the performance of the model before and after fine-tuning. As shown in Fig. 12, we apply Ambiguity Type Classification to both the original Llama3.1-8B (without fine-tuning) and it fine-tuned version with Abg-SciQA. The results indicate that fine-tuning enables the model to cor
794	ectly classify the ambiguity type. Conversely, in Fig. 13, despite undergoing the same fine-tuning
795	process, the model fails to accurately predict the ambiguity type, demonstrating its limitations in
796	JEITAIII CASES.
797	
798	
799	
800	
202	
802	
804	
805	
808	
807	
808	
809	

810 811 812 813 814 815 816	
817	
818	
819	
820	
821	
822	
823	
824	
825	LLM Evaluation
826	Revised Story: {revised story}, Ambiguous Question: {ambiguous question}, Ambiguous answer1: {answer for
827	ambiguous question }, Ambiguous answer2: { answer for ambiguous question }, Clarification answer2: { clarificati answer2: { clarification answer2: { clarification answer2:
828	The ambiguous question has four types: Lexical. Syntactic. Incomplete. and Contextual Ambiguity
829	Example1: Lexical Ambiguity Question: {definition and a question example}
830	clarification question: { clarification question },
831	Example 2: Syntactic Ambiguity Question: {definition and a question example }
832	clarification question: {clarification question},
833	Examples: Incomplete Ambiguity Question: {definition and a question example}
834	Example4: Contextual Ambiguity Question: { definition and a question example }
835	clarification question: { clarification question }.
836	Instructions:
837	Verify if the ambiguous question aligns with the definition of " <lexical, contextual,="" incomplete="" or="" syntactic,=""></lexical,>
838	Ambiguity". If the ambiguous question overlaps with other types of ambiguity, please directly output "False". If the
830	(ambiguous answer 1), `ambiguous answer 2`, `c answer1`, `c answer2`), assess its correctness and provide an
840	explanation with a supporting sentence from the `revised_story`.
841	Output Format:
842	Match with <lexical, contextual,="" incomplete="" or="" syntactic,="">Ambiguity Definition: <True or False></lexical,>
843	Explanation: <reasons align="" aligns="" ambiguity="" does="" not="" or="" question="" the="" why="" with=""> Correctness of ambiguous answer 1: < Correct or Incorrect></reasons>
844	Explanation1: < Reasons why ambiguous_answer_1 is correct or incorrect >
845	Correctness of ambiguous_answer_2: < Correct or Incorrect >
846	Explanation2: < Reasons why ambiguous_answer_2 is correct or incorrect>
847	Correctness of c_answer1: < Correct or Incorrect >
848	Explanation:: < Keasons why c_answer1 is correct or incorrect> Correctness of c_answer2: < Correct or Incorrect>
849	Explanation4: < Reasons why c_answer2 is correct or incorrect >
850	
851	Figure 8: The prompt for Ilm evaluation in Δ ba-SciOA
852	rigure 8. The prompt for thir evaluation in Abg-SerQA.
853	
854	
855	
856	
857	
858	
850	
860	
861	
862	
863	
500	

865 866 867 868 869 870 871 872 873 874 { "id": 464, 875 "story": "Petroleum, consisting of crude oil and natural gas, seems 876 to originate from organic matter in marine sediment. Microscopic 877 organisms settle to the seafloor and accumulate in marine mud. 878 The organic matter may partially decompose, using up the dissolved oxygen in the sediment. As soon as the oxygen is gone, 879 decay stops and the remaining organic matter is preserved. 880 Continued sedimentation-the process of deposits' settling on the 881 sea bottom-buries the organic matter and subjects it to higher 882 temperatures and pressures.", 883 "target_turn": { "question": "What happens to the organic matter in marine sediment 884 over time?", 885 }, 886 "ambiguity_turn": { 887 "ambiguity": "ambiguous", 888 "ambiguity_type": "Lexical Ambiguity" 889 }, "clarification_turn": { 890 "question": "Are you asking about the initial decomposition process 891 of the organic matter or its transformation into oil and gas?" 892 893 "answers": [894 "clr_ans": "The initial decomposition process of the organic 895 matter in marine sediment.", 896 "org_ans": "The organic matter may partially decompose, using up 897 the dissolved oxygen in the sediment and accumulate in marine 898 mud." }, 899 { 900 "clr_ans": "The transformation of the organic matter into oil and 901 gas.", 902 "org_ans": "The organic matter is subjected to higher 903 temperatures and pressures, which convert it to oil and gas." 904 }] 905 } 906 } 907 908 Figure 9: An example of Abg-SciQA data structure 909 910 911 912 913 914 915 916 917

Ambiguity Detection on Abg-SciQA Good Case

Story:

924 925

By 1850, the United States possessed roughly 9,000 miles of railroad track; ten years later, it 926 927 had over 30,000 miles, more than the rest of the world combined. Much of the new construction during the 1850s occurred west of the Appalachian Mountains—over 2,000 miles 928 929 in the states of Ohio and Illinois alone. The effect of the new railroad lines rippled outward 930 through the economy, fundamentally transforming both trade routes and agricultural practices.\n\n\n\n The new railroad networks shifted trade dynamics by redirecting western 931 932 trade from the south to the east, significantly impacting the economic relationships of the time. In 1840, most northwestern grain was shipped south down the Mississippi River to the 933 bustling port of New Orleans. However, this route was fraught with difficulties: low water 934 levels made steamboat travel hazardous in the summer, and ice shut down traffic in winter. 935 Products such as lard, tallow, and cheese spoiled quickly in New Orleans' hot and humid 936 warehouses. Increasingly, traffic from the Midwest flowed from west to east, over the new rail 937 lines. This shift moved trade away from New Orleans and made Chicago a crucial trade hub, 938 linking the farms of the upper Midwest to New York and other eastern cities by more than 939 2,000 miles of track in 1855. As a result, while the value of goods shipped by river to New 940 Orleans continued to increase, the South's overall share of western trade dropped dramatically. 941 Simultaneously, the new rail networks empowered farmers along the tracks to specialize in 942 crops that they could market in distant locations, enhancing the commercial orientation of 943 agriculture and shifting the overall economic landscape. For instance, before the railroad 944 reached Tennessee, the state produced about 25,000 bushels (or 640 tons) of wheat, which sold 945 for less than 50 cents a bushel. Once the railroad came, farmers in the same counties grew 946 400,000 bushels (over 10,000 tons) and sold their crop at a dollar a bushel. A sharp rise in 947 demand for grain abroad also encouraged farmers in the Northeast and Midwest to become 948 more commercially oriented. Wheat, which in 1845 commanded \$1.08 a bushel in New York 949 City, fetched \$2.6 in 1855; similarly, the price of corn nearly doubled. Farmers responded by 950 specializing in cash crops, borrowing to purchase more land, and investing in equipment to 951 increase productivity.... 952 **Target Ouestion:** 953 How did the new rail networks transform the trade dynamics in the western United States 954 during the 1850s? 955 **Detected Ambiguity After Fine-Tuning:** 956 Unambiguous 957 **Explanation:** 958 The question is clear and unambiguous in its wording. It specifies the subject (new rail 959 networks), the geographical area (western United States), and the time period (1850s). 960 **Detected Ambiguity After Fine-Tuning:** 961 Ambiguous 962 Explanation: 963 The question is ambiguous because there are two possible answers for the target question. 964 **Correct Ambiguity:** 965 Ambiguous 966

Figure 10: The good case for Ambiguity Detection on Abg-SciQA

- 970
- 971

967 968

Ambiguity Detection on Abg-SciQA Bad Case

Story:

977 978

Because of industrialization, but also because of a vast increase in agricultural output without 979 which industrialization would have been impossible, Western Europeans by the latter half of 980 the nineteenth century enjoyed higher standards of living and longer, healthier lives than most 981 of the world's peoples. In Europe as a whole, the population rose from 188 million in 1800 to 982 400 million in 1900. By 1900, virtually every area of Europe had contributed to the 983 tremendous surge of population, but each major region was at a different stage of demographic 984 change.Improvements in the food supply continued trends that had started in the late 985 seventeenth century. New lands were put under cultivation, while the use of crops of American 986 origin, particularly the potato, continued to expand. Setbacks did occur. Regional agricultural 987 failures were the most common cause of economic recessions until 1850, and they could lead 988 to localized famine as well. A major potato blight (disease) in 1846-1847 led to the deaths of at 989 least one million persons in Ireland and the emigration of another million, and Ireland never 990 recovered the population levels the potato had sustained to that point. Bad grain harvests at the 991 same time led to increased hardship throughout much of Europe. After 1850, however, the 992 expansion of foods more regularly kept pace with population growth, though the poorer 993 classes remained malnourished. Two developments were crucial. First, the application of 994 science and new technology to agriculture increased. Led by German universities, increasing 995 research was devoted to improving seeds, developing chemical fertilizers, and advancing 996 livestock. After 1861, with the development of land-grant universities in the United States that 997 had huge agricultural programs, American crop-production research added to this mix. 998 Mechanization included the use of horse-drawn harvesters and seed drills, many developed 999 initially in the United States. It also included mechanical cream separators and other food-1000 processing devices that improved supply. The second development involved industrially based 1001 transportation. With trains and steam shipping, it became possible to move foods to needy 1002 regions within Western Europe quickly. Famine (as opposed to malnutrition) became a thing of 1003 the past. Many Western European countries, headed by Britain, began also to import increasing 1004 amounts of food, not only from Eastern Europe, a traditional source, but also from the Americas, Australia, and New Zealand. Steam shipping, which improved speed and capacity, as well as new procedures for canning and refrigerating foods (particularly after 1870), was 1007 fundamental to these developments. ... 1008 **Target Question:** 1009 What caused the food supply to increase in most of Western Europe during the nineteenth 1010 century? 1011 **Detected Ambiguity After Fine-Tuning:** 1012 Unambiguous 1013 **Explanation:**

- 1014 The question is clear and unambiguous in its wording.
- 1015 Detected Ambiguity After Fine-Tuning:
- 1016 Unambiguous
- 1017 Explanation:
- The question is clear and not ambiguous.
- 1019 Correct Ambiguity:
- Ambiguous(There are two possible answers)
- 1023 1024
- Figure 11: The bad case for Ambiguity Detection on Abg-SciQA
- 1025

Ambiguity Type Classification on Abg-SciQA Good Case

1031 Story:

1030

1032 Phytoplankton are minute, free-floating aquatic plants. In addition to the marked changes in 1033 abundance observed in phytoplankton over the course of a year, there is also a notable change 1034 in species composition. This seasonal shift in the dominant species is known as seasonal 1035 succession and occurs in diverse locations. Under seasonal succession, certain species 1036 dominate for a period before being replaced by others, a cycle that repeats annually. ... These organic metabolites likely include various classes of organic compounds. Some of these are toxins, such as those released by dinoflagellates (a type of plankton) during red tides, which 1039 inhibit the growth of other photosynthetic organisms. During red tides, dinoflagellate populations grow so extensively that the water takes on a reddish-brown hue due to the sheer 1040 number of cells. While each cell secretes only a small amount of toxin, the high population 1041 density results in toxin concentrations that can be lethal to many marine creatures. This toxin 1043 can also accumulate in filter-feeding organisms like clams and mussels, making them toxic to humans. Importantly, no special mechanism is required for this population explosion; the 1045 sheer abundance of dinoflagellates alone is sufficient to alter water coloration. Another class of metabolite is vitamins. Certain phytoplankton species require specific vitamins for survival, 1046 with considerable variation in these requirements across species. The B vitamins—especially 1047 B12, thiamine, and biotin—are among the most essential. Some species cannot thrive unless 1048 specific vitamins are available in the water, which are produced by other species. This 1049 dependency can lead to a succession pattern where a vitamin-producing species appears first, 1050 followed by vitamin-dependent species. Other organic compounds, such as amino acids, 1051 carbohydrates, and fatty acids, may also influence phytoplankton succession. Although 1052 laboratory studies show that different phytoplankton species vary in their ability to produce 1053 and require essential vitamins, the extent of this influence in natural marine ecosystems 1054 remains unclear. Additionally, selective grazing by herbivorous crustaceans, such as copepods 1055 and invertebrate larvae, can affect species composition by selectively feeding on certain 1056 phytoplankton. Increasing evidence suggests that multiple factors interact simultaneously to 1057 drive species succession, with the relative importance of each factor varying based on species 1058 and environmental conditions. **Target Question:** How does the population size of dinoflagellates during red tides affect the environment according to the passage? 1062 **Ambiguity Type Classification Before Fine-Tuning:** 1063 Lexical Ambiguity 1064 **Explanation:** The question is ambiguous because the term "population" can refer to either the number of dinoflagellates or their environmental impact, leading to multiple possible interpretations. 1067 **Ambiguity Type Classification After Fine-Tuning:** 1068 Contextual Ambiguity 1069

1070 Explanation:

1079

The ambiguity arises because the question lacks specificity regarding whether it is asking about the direct ecological effects of the dinoflagellate population or the indirect consequences mediated by toxins and food chain disruptions. Different interpretations could lead to different answers.

 Correct Ambiguity Type:

 1075

 1076

 1077

 1078

 Figure 12: The good case for Ambiguity Type Classification on Abg-SciQA

1080 1081 1082 1083 1084 1087 1088 Ambiguity Type Classification on Abg-SciQA Bad Case 1089 **Story:** 1090 Earth has abundant water in its oceans but very little carbon dioxide in its relatively thin 1091 atmosphere. By contrast, Venus is very dry and its thick atmosphere is mostly carbon dioxide. The original atmospheres of both Venus and Earth were derived at least in part from gases 1093 spewed forth, or outgassed, by volcanoes. The gases that emanate from present-day volcanoes 1094 on Earth, such as Mount Saint Helens, are predominantly water vapor, carbon dioxide, and 1095 sulfur dioxide. These gases suggest the possible original atmospheric compositions of both Venus and Earth, emphasizing why Venus now has a thick carbon dioxide atmosphere while Earth does not. Much of the water on both planets is also thought to have come from impacts 1098 from comets, icy bodies formed in the outer solar system.\n\n\nln fact, water probably once 1099 dominated the Venusian atmosphere. Venus and Earth are similar in size and mass, so 1100 Venusian volcanoes may well have outgassed as much water vapor as on Earth, and both 1101 planets would have had about the same number of comets strike their surfaces. Studies of how stars evolve suggest that the early Sun was only about 70 percent as luminous as it is now, so 1102 the temperature in Venus' early atmosphere must have been quite a bit lower. Thus water vapor 1103 would have been able to liquefy and form oceans on Venus. But if water vapor and carbon 1104 dioxide were once so common in the atmospheres of both Earth and Venus, what became of 1105 Earth's carbon dioxide? And what happened to the water on Venus?... But Venus being closer 1106 to the Sun than Earth is, enough of the liquid water on Venus would have vaporized to create a 1107 thick cover of water vapor clouds. Since water vapor is a greenhouse gas, this humid 1108 atmosphere, perhaps denser than Earth's present-day atmosphere, would have efficiently 1109 trapped heat from the Sun. At first, this would have had little effect on the oceans of Venus... 1110 Over time, the rising temperatures would have leveled off, solar ultraviolet radiation having 1111 broken down atmospheric water vapor molecules into hydrogen and oxygen. With all the 1112 water vapor gone, the greenhouse effect would no longer have accelerated.\n\n 1113 **Target Question:** 1114 What evidence from the passage suggests that Venus may have once had substantial water? 1115 **Ambiguity Type Classification Before Fine-Tuning:** 1116 Lexical Ambiguity 1117 **Explanation:** 1118 The target question is ambiguous because the word "substantial" has multiple meanings. 1119 However, there is no keyword 'substantial' in the context. 1120 **Ambiguity Type Classification After Fine-Tuning:** 1121 Incomplete Ambiguity 1122 **Explanation:** 1123 The target question lacks essential contextual information, so it's incomplete ambiguity. 1124 **Correct Ambiguity Type:** 1125 Contextual Ambiguity 1126 1127 Figure 13: The bad case for Ambiguity Type Classification on Abg-SciQA 1128 1129 1130 1131 1132