

# Aligning Language Models to Explicitly Handle Ambiguity

Anonymous ACL submission

## Abstract

In spoken languages, utterances are often shaped to be incomplete or vague for efficiency. This can lead to varying interpretations of the same input, based on different assumptions about the context. To ensure reliable user-model interactions in such scenarios, it is crucial for models to adeptly handle the inherent ambiguity in user queries. However, conversational agents built upon even the most recent large language models (LLMs) face challenges in processing ambiguous inputs, primarily due to the following two hurdles: (1) LLMs are not directly trained to handle inputs that are too ambiguous to be properly managed; (2) the degree of ambiguity in an input can vary according to the intrinsic knowledge of the LLMs, which is difficult to investigate. To address these issues, this paper proposes a method to align LLMs to explicitly handle ambiguous inputs. Specifically, we introduce a proxy task that guides LLMs to utilize their intrinsic knowledge to self-disambiguate a given input. We quantify the information gain from the disambiguation procedure as a measure of the extent to which the models perceive their inputs as ambiguous. This measure serves as a cue for selecting samples deemed ambiguous from the models' perspectives, which are then utilized for alignment. Experimental results from several question-answering datasets demonstrate that the LLMs fine-tuned with our approach are capable of handling ambiguous inputs while still performing competitively on clear questions within the task.

## 1 Introduction

Large Language Models (LLMs) (Ouyang et al., 2022; Team et al., 2023; Achiam et al., 2023) have demonstrated remarkable capabilities in text generation, proving particularly effective for question-answering (QA) tasks (Zhang et al., 2023; Etezadi and Shamsfard, 2023). QA systems in the wild are frequently confronted with unexpected inputs from

When was the last time UGA won a national championship?

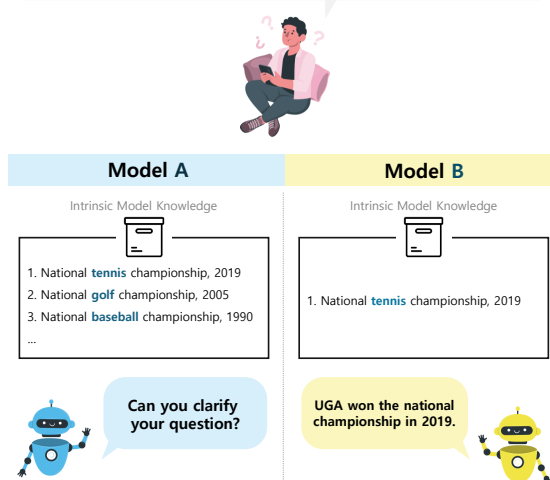


Figure 1: An example of an ambiguous query from AmbigQA (Min et al., 2020). The phrase "national championship" poses diverse denotations, causing ambiguity. Model possessing various related knowledge of the query could perceive it as ambiguous (left). On the other hand, if the model does not have sufficient related knowledge (right), the query can be perceived as unambiguous. Thus, the degree of ambiguity perceived by the model may vary even with identical inputs.

users, such as unanswerable (Kim et al., 2023b; Yin et al., 2023) or ambiguous questions (Cole et al., 2023; Lee et al., 2023; Kim et al., 2023a). To build a reliable user-friendly model, it is essential for the model to robustly handle such inputs. In this work, we seek to extend the scope of research to effectively handle invalid inputs. Specifically, we focus on managing "ambiguity" (Gleason, 1963; Mackay and Bever, 1967), which poses a significant challenge in Natural Language Processing (NLP) (Jurafsky, 1996).

Ambiguity refers to cases where an expression conveys multiple denotations (Wasow et al., 2005). Users may pose queries with clear intentions that, possibly due to insufficient domain knowledge,

058 result in ambiguous requests. If the model arbitrary-ly responds to such ambiguity, there is a risk  
059 of misinterpreting the user’s original intent, po-  
060 tentially harming the model’s reliability. This is  
061 especially pronounced in sensitive domains such  
062 as legal (Schane, 2002; Choi, 2024) or medical  
063 (Stevenson and Guo, 2010; Gyori et al., 2022) do-  
064 mains, where misinterpretations can lead to serious  
065 drawbacks. Despite the importance, approaches to  
066 robustly manage ambiguity are still significantly  
067 unexplored. In this paper, we endeavor to utilize  
068 the model’s intrinsic knowledge to align the model  
069 in a manner that effectively handles ambiguity.  
070

071 Properly processing ambiguous inputs is chal-  
072 lenging primarily due to the following two hurdles.  
073 Firstly, models are **not directly trained to explic-  
074 itly express ambiguity**. Even if a model perceives  
075 ambiguity, it is challenging to verify the recogni-  
076 tion without explicit feedback. The second chal-  
077 lenge is that the degree of ambiguity for the query  
078 can **vary depending on the intrinsic knowledge**  
079 of the model. Consider the scenario depicted in Fig-  
080 ure 1. The initial query is ambiguous as the phrase  
081 "national championship" poses various denotations,  
082 such as "national tennis championship" or "national  
083 golf championship". If a model possesses compre-  
084 hensive knowledge across the possible denotations,  
085 it is plausible for the model to recognize the ambi-  
086 guity (left). However, if the model’s knowledge is  
087 limited to "national tennis championship", it would  
088 perceive the query as unambiguous (right). There-  
089 fore, it is essential to verify whether the input is  
090 deemed ambiguous from the model’s point of view.

091 To overcome these issues, this paper proposes  
092 a method to align models to explicitly handle am-  
093 biguous queries. Specifically, we design a proxy  
094 task that guides the model to self-disambiguate a  
095 given query by utilizing its intrinsic knowledge.  
096 Then, we quantify the information gain from the  
097 disambiguation as an implicit measure of the extent  
098 to which the models perceive their inputs as am-  
099 biguous. This measure serves as a cue for selecting  
100 samples deemed ambiguous from the model’s per-  
101 spective, which are then utilized for alignment. Ex-  
102 perimental results from several QA datasets demon-  
103 strate that the alignment process enables the model  
104 to properly clarify ambiguous inputs while main-  
105 taining its inherent capabilities. The findings un-  
106 derscore the value of assessing the perceived ambi-  
107 guity, rather than relying solely on the ground-truth  
108 ambiguity. Furthermore, to provide a comprehen-  
109 sive framework for assessing ambiguity, we con-

110 struct a new dataset dubbed AmbigTriviaQA. The  
111 dataset facilitates a more extensive evaluation of  
112 models’ robustness in addressing ambiguity, thus  
113 contributing to the further expansion of related re-  
114 search.

## 2 Related Work 115

**Ambiguity in NLP** An expression is defined as  
116 ambiguous if it has two or more distinct denotations  
117 (Wasow et al., 2005). Ambiguity challenges NLP  
118 applications by obscuring the intended meaning  
119 of expressions, leading to difficulties in accurately  
120 performing specific tasks. Efforts addressing this  
121 issue span across various domains, including ma-  
122 chine translation (Pilault et al., 2023), coreference  
123 resolution (Poesio and Artstein, 2005; Yuan et al.,  
124 2023), and natural language inference (Liu et al.,  
125 2023).  
126

127 The challenge intensifies in the scope of QA as  
128 ambiguous questions may yield various answers,  
129 potentially not aligning with the user’s initial intent.  
130 Min et al. (2020) introduce the AmbigQA dataset to  
131 tackle ambiguity in open-domain QA and Stelmakh  
132 et al. (2022) expands it to long-form generation.  
133 Furthermore, Cole et al. (2023) discovered that  
134 quantifying sampling repetition presents a reliable  
135 uncertainty measure for ambiguity, while Kim et al.  
136 (2023a) generates tree-of-clarification (ToC) that  
137 refines ambiguity within the inputs. As we share  
138 the goal of handling ambiguity, we adopt a novel  
139 approach of directly aligning the model to address  
140 ambiguity.

**Alignment of LLMs** LLMs are fundamentally  
141 trained through causal language modeling, a pro-  
142 cess essential for understanding and generating text  
143 of high fluency and consistency. To better harness  
144 these models, approaches have been developed to  
145 align them with human preferences (Leike et al.,  
146 2018; Ji et al., 2023b). This has taken various  
147 forms, notably through Reinforcement Learning  
148 from Human Feedback (RLHF) (Ouyang et al.,  
149 2022; Bai et al., 2022a; Chakraborty et al., 2024),  
150 as well as Supervised Fine-tuning (SFT) (Dong  
151 et al., 2023; Yang et al., 2023; Zhou et al., 2024).  
152

153 Previous works focused on preferences such as  
154 helpfulness (Ding et al., 2023; Köpf et al., 2023; Xu  
155 et al., 2024) and safety (Bai et al., 2022b; Ji et al.,  
156 2023a; Liu et al., 2024b). Recent studies have con-  
157 centrated on the factuality (Yang et al., 2023; Tian  
158 et al., 2024), avoiding hallucinations. Building  
159 on this foundation, our research extends the scope

and aims to align models to effectively understand and handle ambiguities, which is a relatively unexplored area within the field of model alignment. This stands in contrast to previous methods which typically bypass the nuanced interpretation of ambiguous contexts inherent in language.

**Data Quality Control for Alignment** Data-centric AI (Chu et al., 2016; Majeed and Hwang, 2023; Kumar et al., 2024) emphasizes the importance of data quality in training models. In the context of the instruction following, LIMA (Zhou et al., 2024) demonstrates that models can be effectively aligned even with 1,000 human-curated samples. Similarly, AlpaGasus (Chen et al., 2024) shows effective alignment can be achieved by utilizing only a small subset of the Alpaca dataset (Taori et al., 2023) selected by ChatGPT. Various approaches for data selection have been explored, including those based on pre-defined quality factors such as length and complexity (Liu et al., 2024a), and utilizing gradient similarity from validation sets as a selection criterion (Xia et al., 2024). In this paper, we distinctly define data quality to effectively align models to handle ambiguity. To do so, we make use of the model’s perceived ambiguity as an implicit cue for data quality.

### 3 Methodology

The objective of our approach is to align models to explicitly handle potentially ambiguous inputs leveraging intrinsic model knowledge. To this end, we propose a four-stage alignment pipeline, depicted in Figure 2. In this section, we first formulate the problem and describe each stage in detail.

**Problem Formulation** The goal of a QA task is to generate a factually correct answer  $y$ , given an unambiguous input  $x_{\text{unambig}}$ , a pre-defined inference template  $t(\cdot)$ , and a language model  $M$ . As we expand our input scope to ambiguous queries  $x_{\text{ambig}}$ , the model is expected to generate a clarification request  $y_{\text{clarify}}$ <sup>1</sup> for  $x_{\text{ambig}}$  to resolve the ambiguity, where the user is best positioned to clarify their intent.

<sup>1</sup>We have considered various approaches to handle ambiguity but were concluded to be impractical. Arbitrarily offering one of the valid answers may fail to reflect the user’s intent, and presenting all possible answers is often impractical due to the potentially vast number of valid answers.

#### 3.1 Explicit Prediction Stage

This initial stage involves assessing whether the model can appropriately handle each sample and identifying samples that the model currently fails to explicitly manage. By comparing the model’s prediction with the ground-truth label, samples are categorized based on their response accuracy. We collect  $n$  correct samples, which the model can properly handle, as  $D_{\text{correct}} = \{(x_{\text{correct}}^i, y_{\text{correct}}^i)\}_{i=1}^n$  and incorrect samples are classified as  $D_{\text{incorrect}}$ .

#### 3.2 Implicit Ambiguity Detection Stage

The objective of this stage is to identify samples that the model perceives as ambiguous from  $D_{\text{incorrect}}$ . Given that it is challenging for the model to explicitly express ambiguity, we construct a proxy task to estimate the ambiguity from the model’s point of view.

The proxy task is designed to self-disambiguate  $x$  and implicitly measure the perceived ambiguity. Specifically, the model is first prompted to generate a disambiguation  $\hat{x}_{\text{disambig}}$  for the input  $x$ . In this process, the model leverages its intrinsic knowledge related to  $x$  and generates further details. If  $x$  lacks specifications and the model possesses related knowledge necessary to compensate, then  $\hat{x}_{\text{disambig}}$  would yield a higher certainty (lower entropy) for the model. On the other hand, if  $x$  requires no specification or the model lacks the necessary knowledge,  $\hat{x}_{\text{disambig}}$  would exhibit a similar level of uncertainty to  $x$ . To quantify the uncertainty associated with  $x$  and  $\hat{x}_{\text{disambig}}$ , we employ the model’s average entropy (Malinin and Gales, 2021; Abdar et al., 2021). Formally, the entropy of an output distribution is defined as follows:

$$\mathcal{H}_{x,i} = - \sum_{v \in \mathcal{V}} p_{x,i}(v) \log p_{x,i}(v) \quad (1)$$

where  $p_{x,i}(v)$  is the probability of the  $i^{\text{th}}$  token  $v$  of a sentence  $x$  from the full vocabulary set  $\mathcal{V}$ . The average entropy for  $x$  can be defined as:

$$\mathcal{H}_x = \frac{1}{I} \sum_i \mathcal{H}_{x,i} \quad (2)$$

where  $x$  is composed of  $I$ -tokens. We quantify the changes in input uncertainty by the difference in average entropy, which we define as **information gain** (INFOGAIN). The INFOGAIN from the disambiguation can be defined as the following:

$$\text{INFOGAIN}_{x, \hat{x}_{\text{disambig}}} = \mathcal{H}_x - \mathcal{H}_{\hat{x}_{\text{disambig}}} \quad (3)$$

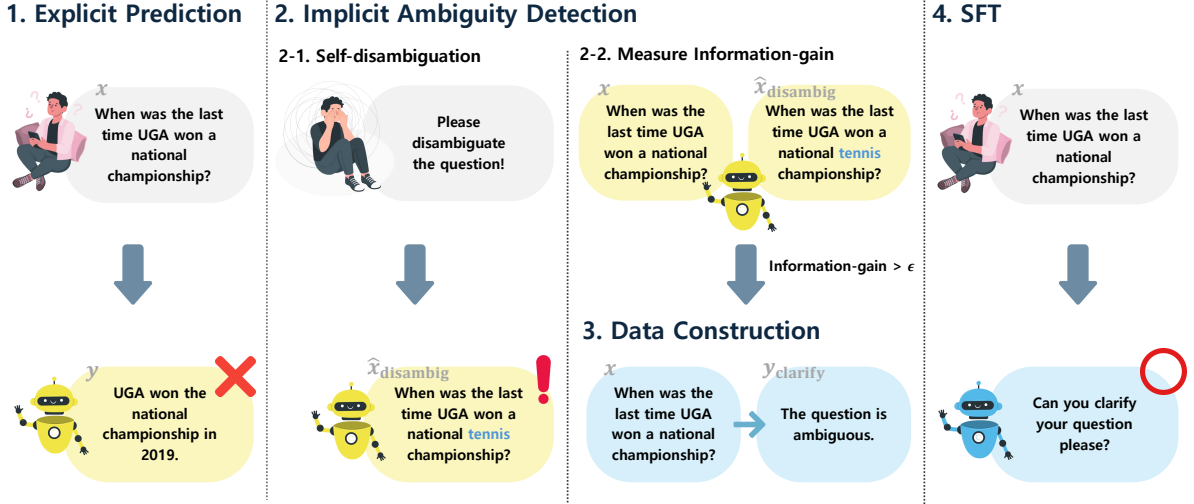


Figure 2: The overall process of our four-stage alignment pipeline. We initially filter samples that the model cannot explicitly handle (Stage 1), and self-disambiguate them to measure the information gain (Stage 2). Samples with high information gain are deemed ambiguous and utilized for supervised fine-tuning (Stage 3 & 4).

If the disambiguation results in a meaningful specification by utilizing intrinsic knowledge, a substantial INFOGAIN would be measured, suggesting that the model considers  $x$  as ambiguous. Conversely, a negligible INFOGAIN indicates that the model does not perceive  $x$  as ambiguous. Samples with INFOGAIN greater than the threshold  $\epsilon$  are classified as ambiguous, denoted as  $x_{\text{ambig}}$ .

### 3.3 Data Construction Stage

In this stage, we construct datasets for the alignment process. This involves labeling  $m$  samples identified as ambiguous and constructing an ambiguous dataset  $D_{\text{ambig}} = \{(x_{\text{ambig}}^j, y_{\text{clarify}})\}_{j=1}^m$ .  $y_{\text{clarify}}$  serves as the ground-truth label for ambiguous samples, which are randomly selected from pre-defined clarification requests stipulated in Appendix D. To prevent the potential loss of the model’s existing knowledge, we also incorporate  $D_{\text{correct}}$  for training. We balance the number of samples from both datasets so that  $n = m$ . The final training dataset is thus established as  $D = D_{\text{correct}} + D_{\text{ambig}}$ .

### 3.4 Supervised Fine-tuning (SFT) Stage

Utilizing the dataset  $D = \{(x^k, y^k)\}_{k=1}^{n+m}$ , the model is trained to generate ground-truth label  $y$  for input  $x$ , employing the identical inference template  $t(\cdot)$ . The model  $M$  with parameter  $\theta$  is trained as follows:

$$\min_{\theta} \sum_{(x,y) \in D} \sum_{i=1}^{|y|} -\log M_{\theta}(y_i | y_{<i}, t(x)) \quad (4)$$

## 4 Experimental Setting

### 4.1 Datasets

The capability of the model to perform within the trained domain is pivotal. However, its ability to generalize to out-of-distribution (OOD) is essential for real-world applicability, as queries that deviate from the training data are frequently confronted in the wild. To this end, we employ one training dataset and three OOD test sets to evaluate in diverse domains. All the datasets include both ambiguous and unambiguous queries.

**AmbigQA** Introduced by Min et al. (2020), AmbigQA is a derivative of the Natural Questions dataset (Kwiatkowski et al., 2019), designed to verify data points deemed ambiguous. The dataset covers diverse sources of ambiguity such as event and entity references. We set AmbigQA as the in-domain dataset and utilize it for training.

**SituatedQA** SituatedQA (Zhang and Choi, 2021) specifically focuses on temporal and geographic ambiguity from the input query. As the cause of ambiguity and its construction process are distinct, we assess performance on the temporal split and the geographic split separately, denoted as Temp and Geo, respectively.

**AmbigTriviaQA** Since there are limited datasets for evaluating ambiguity in open-domain QA, we construct a new dataset, namely AmbigTriviaQA. By taking questions from the widely-used TriviaQA dataset (Joshi et al., 2017), we prompt



306	gpt-3.5-turbo to ambiguat the initial query and	354
307	verify the results. More details on dataset construc-	355
308	tion are described in Appendix B.	356
309	<b>4.2 Baselines</b>	357
310	To assess the effectiveness of our approach, we es-	358
311	tablish two sets of baselines: inference-only meth-	359
312	ods and trained methods. Further implementation	360
313	details are described in Appendix A.	361
314	<b>Inference-Only Methods</b> Inference-only meth-	362
315	ods address ambiguity by directly prompting the	363
316	model. We employ naïve prompting (NAÏVE) as	364
317	a fundamental baseline, applying a simple QA	365
318	prompt. Furthermore, we explore ambiguity-aware	366
319	prompting (AMBIGUITY-AWARE), which addition-	367
320	ally provides instructions on handling ambiguity.	368
321	We also examine SAMPLE REPETITION (Cole	369
322	et al., 2023) by measuring the consistency of the	370
323	sampled generations. Finally, we evaluate SELF-	371
324	ASK (Amayuelas et al., 2023), where the model	372
325	initially generates an answer and subsequently de-	373
326	termines the ambiguity based on the generation.	374
327	<b>Trained Methods</b> Given the lack of directly com-	375
328	parable prior work, we compare a fine-tuned base-	376
329	line wherein the model is trained with the in-	377
330	domain training set. FULL-SET applies the full	378
331	in-domain training dataset and SUBSET is trained	379
332	on a randomly selected subset of size $n + m$ , which	380
333	is equivalent in size to our approach. Additionally,	381
334	we compare HONESTY-TUNED (Yang et al., 2023),	382
335	which takes a similar approach utilizing an implicit	383
336	measure named "expected accuracy" to estimate	384
337	the model's factuality. The expected accuracy is	385
338	measured as the average accuracy of sampled pre-	386
339	diction for a single input query. We have revised	387
340	the method so that incorrect samples based on "ex-	388
341	pected accuracy" are selected from the ground-truth	389
342	ambiguous samples. Although HONESTY-TUNED	390
343	shares similarities with our approach in the way	391
344	of estimating the model's knowledge with an im-	392
345	PLICIT measure, a notable distinction lies in our main	393
346	focus on handling ambiguity beyond factuality.	394
347	<b>4.3 Evaluation Metrics</b>	395
348	As we expand the input scope to possibly ambigu-	396
349	ous questions, the model should be capable of han-	397
350	dling both unambiguous and ambiguous queries	398
351	simultaneously. Therefore, we employ two widely	399
352	used evaluation metrics to assess performance on	
353	both types of inputs. A successful alignment should	
	preserve the model's capability to handle unam-	
	biguous inputs while successfully managing am-	
	biguous queries. All evaluations are conducted	
	by comparing the greedy generation to the ground	
	truth.	
	<b>Unambiguous Accuracy (Unambig. Acc.)</b>	
	While expanding the task scope, it remains cru-	
	cial for the model to preserve the ability to handle	
	unambiguous inputs. Thus, our analysis persists	
	in exclusively evaluating the model's accuracy in	
	processing unambiguous queries. We measure the	
	quality of the generation by employing RougeL <sup>2</sup>	
	(Lin and Och, 2004) with all the possible answers,	
	where the prediction is regarded as correct if the	
	score is above 0.3.	
	<b>Ambiguity Detection F1-score (Ambig. F1)</b>	
	The model should be capable of detecting ambi-	
	guity and generating clarification requests for am-	
	biguous inputs. However, especially for trained	
	methods, models may exhibit biased predictions	
	toward clarification requests. Taking these aspects	
	into account, we evaluate the model's ambiguity	
	detection capability with F1-score, which captures	
	both the precision and recall of prediction, offering	
	a balanced view of the model's ambiguity detec-	
	tion performance. Further details on the detection	
	process are described in Appendix C.	
	<b>4.4 Implementation Details</b>	
	For our experiments, we utilize LLAMA2 7B &	
	13B (Touvron et al., 2023) and MISTRAL 7B	
	(Jiang et al., 2023). We utilized QLoRA (Det-	
	tmers et al., 2023) to facilitate efficient training. Im-	
	plementation details are stipulated in Appendix D.	
	<b>5 Experimental Results</b>	
	The main results of our experiments are presented	
	in Table 1. <b>Inference-only methods exhibit a</b>	
	<b>pronounced deficiency in handling ambiguous</b>	
	<b>queries.</b> Specifically, NAÏVE establish poor per-	
	formance in responding to ambiguous queries, re-	
	sulting in a notably low F1-score. AMBIGUITY-	
	AWARE demonstrates a strong bias towards clari-	
	fication requests, as it achieves a relatively high	
	F1-score at the expense of accuracy. Similarly,	
	SAMPLE REPETITION exhibits a substantial trade-	
	off between F1-score and accuracy. SELF-ASK	
	displays a subpar F1-score, indicating that it is	

<sup>2</sup><https://huggingface.co/spaces/evaluate-metric/rouge>

Method	# Train Samples	AmbigQA		SituatdQA (Geo)		SituatdQA (Temp)		Ambig-TriviaQA	
		Unambig. Acc.	Ambig. F1	Unambig. Acc.	Ambig. F1	Unambig. Acc.	Ambig. F1	Unambig. Acc.	Ambig. F1
<b>LLAMA2 7B</b>									
NAÏVE	0	28.43	0.00	22.53	0.00	<u>21.72</u>	0.00	<b>62.46</b>	0.00
AMBIGUITY-AWARE	0	4.94	<u>68.95</u>	3.95	32.44	1.72	35.53	22.63	61.03
SAMPLE REPETITION	0	5.42	<b>74.96</b>	3.56	34.43	4.40	<u>38.43</u>	29.58	<b>65.88</b>
SELF-ASK	0	26.75	13.41	21.54	8.18	19.68	18.48	<u>61.29</u>	3.84
HONESTY-TUNED	3,088	18.19	68.00	9.49	36.56	7.98	37.57	52.42	54.60
SUBSET	3,088	<u>29.16</u>	64.23	23.72	36.85	13.95	36.74	51.71	55.70
FULL-SET	10,036	<b>37.11</b>	59.98	<b>26.09</b>	<u>41.15</u>	18.93	35.38	58.25	49.96
OURS	3,088	27.23	63.69	<u>24.51</u>	<b>42.05</b>	<b>21.90</b>	<b>40.77</b>	53.41	<u>61.34</u>
<b>MISTRAL 7B</b>									
NAÏVE	0	13.01	0.00	7.51	0.00	10.91	0.00	37.52	0.00
AMBIGUITY-AWARE	0	9.76	54.74	2.17	26.01	5.33	22.48	27.87	44.94
SAMPLE REPETITION	0	4.70	51.21	2.57	29.25	2.15	29.64	25.54	32.54
SELF-ASK	0	13.01	0.00	7.51	0.00	10.91	0.00	37.52	0.00
HONESTY-TUNED	1,382	16.51	<b>69.28</b>	11.66	33.84	8.09	<u>39.16</u>	41.70	<b>65.19</b>
SUBSET	1,382	33.49	60.91	<u>29.05</u>	35.64	19.11	37.38	<u>58.87</u>	54.61
FULL-SET	10,036	<b>43.73</b>	<u>62.85</u>	24.11	<u>40.81</u>	<u>26.76</u>	24.14	<b>65.04</b>	49.63
OURS	1,382	<u>37.23</u>	50.31	<b>32.21</b>	<b>42.18</b>	<b>35.74</b>	<b>40.17</b>	58.14	<u>58.93</u>
<b>LLAMA2 13B</b>									
NAÏVE	0	31.33	0.00	<u>22.53</u>	0.00	<u>22.90</u>	0.00	60.49	0.14
AMBIGUITY-AWARE	0	3.37	<u>70.44</u>	3.16	33.10	2.22	36.66	16.96	<b>64.17</b>
SAMPLE REPETITION	0	10.00	<b>71.10</b>	6.32	32.85	9.45	37.87	39.41	<u>63.40</u>
SELF-ASK	0	31.33	0.00	<u>22.53</u>	0.00	<u>22.90</u>	0.00	60.49	0.14
HONESTY-TUNED	3,216	17.83	68.57	3.16	33.29	3.58	38.03	48.13	60.39
SUBSET	3,216	35.18	62.89	23.12	36.19	20.50	<u>39.00</u>	59.72	57.26
FULL-SET	10,036	<b>43.61</b>	62.87	23.12	<u>38.37</u>	19.68	24.07	<b>66.80</b>	48.81
OURS	3,216	<u>37.83</u>	58.15	<b>24.51</b>	<b>41.59</b>	<b>24.36</b>	<b>41.09</b>	<u>63.74</u>	55.23

Table 1: Experimental results of in-domain dataset and three OOD datasets. We report unambiguous accuracy (Unambig. Acc.) and ambiguity detection F1-score (Ambig. F1). For each dataset, the **best method** is highlighted in bold and the second-best method is underlined. Our method demonstrates comparable performance in-domain and outperforms all the baselines in the OOD setting.

challenging to resolve ambiguity by explicitly "self-asking" the model.

**Trained methods exhibit enhanced performance overall compared to inference-only approaches.** HONESTY-TUNED struggles to handle ambiguity, as it also demonstrates biased detection performance. This is likely because the implicit measure from HONESTY-TUNED can be influenced by various factors but not specifically ambiguity. The results underscore the necessity of distinct methods for perceiving ambiguity. Compared to HONESTY-TUNED, SUBSET exhibits relatively balanced performance across both metrics. FULL-SET demonstrates the most superior performance among the baselines, particularly in the in-domain setting, as it has access to the ground-truth ambiguity.

**Our approach yields comparable results in in-domain and demonstrates superior OOD per-**

**formances.** Despite employing identical inference templates as NAÏVE, our method demonstrates equal or improved unambiguous accuracy. This indicates the effectiveness of our alignment in managing ambiguity while preserving the inherent capabilities of the model. We can also observe an improvement in unambiguous accuracy, particularly in SituatedQA splits. It is especially surprising given that our method was trained on  $D_{correct}$ , which the model is already capable of handling. Compared to FULL-SET, we note a slight decline in the in-domain performance, an expected result given that FULL-SET is optimized with ground-truth ambiguity of the in-domain data. However, our method outperforms FULL-SET across OOD datasets in F1-score up to 17 points. This discrepancy underscores the effectiveness of utilizing perceived ambiguity for alignment, facilitating superior generalization and robustness. The efficacy

Method	AmbigQA		SituatdQA (Geo)	
	Unambig. Acc.	Ambig. F1	Unambig. Acc.	Ambig. F1
RANDOM	<u>29.40</u>	<b>64.74</b>	18.38	36.23
IMPLICIT	<b>29.52</b>	62.23	<b>24.90</b>	<u>41.19</u>
OURS	27.23	<u>63.69</u>	<u>24.51</u>	<b>42.05</b>

Method	SituatdQA (Temp)		Ambig-TriviaQA	
	Unambig. Acc.	Ambig. F1	Unambig. Acc.	Ambig. F1
RANDOM	16.99	38.76	<b>56.13</b>	50.63
IMPLICIT	<u>21.36</u>	<u>40.11</u>	<u>55.65</u>	<u>59.16</u>
OURS	<b>21.90</b>	<b>40.77</b>	53.41	<b>61.34</b>

Table 2: Ablation results of ambiguous data selection. **In-domain** dataset is highlighted in gray. **The best method** is in bold, and the second best method is underlined. Results show that perceived ambiguity measured by INFOGAIN is an effective cue for data selection.

of leveraging only the data perceived ambiguous (about 32% in the LLAMA2 family and 13% in MISTRAL) emphasizes the importance of data quality over quantity (Zhou et al., 2024; Chen et al., 2024).

## 6 Ablation Study

### 6.1 Impact of INFOGAIN for Data Selection

For a deeper analysis of the influence of INFOGAIN for data selection within our pipeline, we conduct an ablation study by varying the criteria for selecting ambiguous data. While maintaining the same  $D_{\text{correct}}$  for unambiguous samples, we alter the selection of  $m$  samples labeled as ambiguous. We compare the following data selection strategies:

- **Random Selection (RANDOM)** We randomly select  $m$  ground-truth ambiguous samples, without any consideration of INFOGAIN.
- **Implicit Measure-based Selection (IMPLICIT)** We select top- $m$  samples with the largest INFOGAIN among those that are ground-truth ambiguous. It differs from our approach as our method utilizes samples perceived as ambiguous, allowing the potential inclusion of unambiguous samples.

Table 2 is the ablation results on LLAMA2 7B. For AmbigQA, baselines leveraging ground-truth ambiguity slightly outperform our method, which

Method	VAR ( $\uparrow$ )	MCR ( $\downarrow$ )	OAP ( $\uparrow$ )
<b>AmbigQA</b>			
HONESTY-TUNED	<b>73.81</b>	44.49	20.48
SUBSET	62.12	33.05	20.79
FULL-SET	52.82	<b>18.64</b>	<u>21.48</u>
OURS	<u>65.19</u>	<u>30.51</u>	<b>22.65</b>
<b>SituatdQA (Geo)</b>			
HONESTY-TUNED	<b>93.80</b>	58.77	19.34
SUBSET	63.57	<u>35.96</u>	20.35
FULL-SET	72.09	<u>35.96</u>	<u>23.08</u>
OURS	<u>89.15</u>	<b>28.95</b>	<b>31.67</b>
<b>SituatdQA (Temp)</b>			
HONESTY-TUNED	<b>85.73</b>	56.34	<u>18.71</u>
SUBSET	68.61	54.86	15.48
FULL-SET	58.56	<u>46.95</u>	15.53
OURS	<u>72.95</u>	<b>36.08</b>	<b>23.31</b>
<b>AmbigTriviaQA</b>			
HONESTY-TUNED	42.39	<u>11.64</u>	18.73
SUBSET	<u>48.10</u>	17.57	<u>19.83</u>
FULL-SET	38.93	<b>9.81</b>	17.56
OURS	<b>57.45</b>	16.91	<b>23.87</b>

Table 3: VAR, MCR, and OAP of trained methods. A high OAP is preferred, with a high VAR and a low MCR. **The best performance** is in bold, and the second best performance is underlined. Our approach demonstrates the best OAP across all datasets, with the least trade-off between VAR and MCR.

is a similar tendency from Section 5 where FULL-SET exhibit better in-domain performance. However, across OOD datasets, our approach demonstrates significantly superior performance. Specifically, RANDOM demonstrates a notable drop in F1-score by up to 10 points, illustrating the limitations of simply utilizing the ground-truth ambiguity, which might not align with the model’s perceived ambiguity. Furthermore, IMPLICIT surpasses RANDOM by up to 5 points in F1-score, validating the effectiveness of INFOGAIN as a cue for data selection. Finally, our approach outperforms IMPLICIT across the majority of metrics, even with unambiguous samples selected for training, again highlighting the effectiveness of INFOGAIN as the data selection measure.

### 6.2 Analysis on Sample-level Prediction Change

Our method is designed to align the model to generate clarification requests for ambiguous queries. However, the process may lead to a potential

Model Prediction	Ground Truth	Type	Generated Text
Unambig.	Ambig.	$x$ $\hat{x}_{\text{disambig}}$	Who sings don't mess around with jim? Who sings don't mess around with jim, <b>in the 1960s?</b>
Unambig.	Unambig.	$x$ $\hat{x}_{\text{disambig}}$	Who is winner in bigg boss season 5 kannada? Who is the winner of the fifth season of the kannada <b>version of the indian reality television series bigg boss?</b>
Ambig.	Ambig.	$x$ $\hat{x}_{\text{disambig}}$	How many jury members in a criminal trial? How many jury members are required in a criminal trial <b>in the united states?</b>
Ambig.	Unambig.	$x$ $\hat{x}_{\text{disambig}}$	How many pages in a brave new world? How many pages are in the <b>1932 novel</b> brave new world <b>by aldous huxley?</b>

Table 4: Example of initial query  $x$  and its disambiguation  $\hat{x}_{\text{disambig}}$  generated by LLAMA2 7B. **Additional specification** from the model is in bold. Unambig. and Ambig. refers to Unambiguous and Ambiguous, respectively.

trade-off, where the model erroneously generates clarification requests for unambiguous inputs that were previously well-handled. To assess this balance, we introduce three metrics: **Valid Alignment Rate (VAR)** measures the proportion of ambiguous samples incorrectly handled before alignment that are correctly addressed post-alignment and **Misaligned Clarification Rate (MCR)** measures the rate of correct unambiguous samples before training that erroneously generates clarification requests after alignment. A high VAR is desirable, whereas a low MCR is preferred simultaneously. Inspired by Yang et al. (2023), we additionally define **Overall Alignment Performance (OAP)** that measures the balance between VAR and MCR.

$$\text{OAP} = \frac{\text{VAR} + (1 - \text{MCR})}{2} \quad (5)$$

Table 3 compares the results on LLAMA2 7B. HONESTY-TUNED exhibits high VAR but poor MCR, implying a tendency to misinterpret known knowledge as ambiguous. This aligns with the previous results where HONESTY-TUNED displays biased generation towards clarification requests. On the other hand, FULL-SET and SUBSET demonstrate a good MCR and a relatively low VAR. Our method performs superior OAP, successfully addressing ambiguities (high VAR) while preserving existing capabilities (low MCR).

## 7 Self-disambiguation Case Study

Table 4 demonstrates examples of initial query  $x$  and its disambiguation  $\hat{x}_{\text{disambig}}$  generated by LLAMA2 7B. The first example is when  $x$  is inherently ambiguous, yet the model perceives it as unambiguous. Specifically, the model generates hallucination ("in the 1960s") where the song "don't

mess around with jim" was originally released in 1972. This non-factual generation would not provide any information gain to the model, classifying  $x$  as ambiguous. In such a case,  $x$  should be considered "unknown" with no related knowledge within the model. The second and third examples are correctly classified, as the model properly applies its intrinsic knowledge to perceive ambiguity. Regardless of the quantity of additional context generated, the model is capable of verifying its ambiguity. The last example is a misclassification as ambiguous. Despite disambiguation provides factually correct information ("1932 novel" and "by Aldous Huxley") for "brave new world", we speculate that the misclassification may arise from the existence of various media, such as movies and songs, sharing the title "brave new world", leading to an erroneous integration of knowledge.

## 8 Conclusion

In this paper, we present a novel alignment pipeline designed to enhance the ability of LLMs to address ambiguities within queries, leveraging the model's intrinsic knowledge. Our method employs an implicit measure, dubbed INFOGAIN, to quantify ambiguity as perceived by the model. Through alignment based on the measure, the model learns to explicitly handle ambiguous as well as unambiguous queries. Experimental results demonstrate the effectiveness of our alignment, particularly pronounced in out-of-distribution scenarios. Results indicate the importance of alignment based on the model's perceived ambiguity. Future work may explore the extension of this methodology to broader domains and more complex types of ambiguities, further solidifying the role of LLMs in managing the inherent uncertainty present in NLP tasks.



## 556 Limitations

557 For the experiments, we explore the most widely  
558 used models for evaluation, specifically LLAMA2  
559 and MISTRAL. Despite this, a more comprehensive  
560 evaluation encompassing a broader consideration  
561 of LLMs could have enriched our findings, provid-  
562 ing insights across different architectures and capa-  
563 bilities. Larger models in scale could demonstrate  
564 different tendencies and should be explored for fu-  
565 ture work. Furthermore, our work mainly focuses  
566 on supervised fine-tuning (SFT) as the alignment  
567 method. However, alternative methods, such as  
568 Reinforcement Learning from Human Preference  
569 (RLHF) (Ouyang et al., 2022) or Direct Prefer-  
570 ence Optimization (DPO) (Rafailov et al., 2023)  
571 could offer distinct advantages towards our objec-  
572 tive. Finally, the experiments are mainly focused  
573 on short-form QA tasks. The research scope could  
574 be expanded to long-form generation tasks such as  
575 detailed reasoning.

## 576 References

577 Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana  
578 Rezazadegan, Li Liu, Mohammad Ghavamzadeh,  
579 Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U. Ra-  
580 jendra Acharya, Vladimir Makarencov, and Saeid  
581 Nahavandi. 2021. [A review of uncertainty quantifica-  
582 tion in deep learning: Techniques, applications and  
583 challenges](#). *Information Fusion*, 76:243–297.

584 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
585 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
586 Diogo Almeida, Janko Altschmidt, Sam Altman,  
587 Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#).  
588 *arXiv preprint arXiv:2303.08774*.

589 Alfonso Amayuelas, Liangming Pan, Wenhui Chen, and  
590 William Wang. 2023. [Knowledge of knowledge: Ex-  
591 ploring known-unknowns uncertainty with large lan-  
592 guage models](#). *Preprint*, arXiv:2305.13712.

593 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda  
594 Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
595 Stanislav Fort, Deep Ganguli, Tom Henighan,  
596 Nicholas Joseph, Saurav Kadavath, Jackson Kernion,  
597 Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac  
598 Hatfield-Dodds, Danny Hernandez, Tristan Hume,  
599 Scott Johnston, Shauna Kravec, Liane Lovitt, Neel  
600 Nanda, Catherine Olsson, Dario Amodei, Tom  
601 Brown, Jack Clark, Sam McCandlish, Chris Olah,  
602 Ben Mann, and Jared Kaplan. 2022a. [Training  
603 a helpful and harmless assistant with reinforce-  
604 ment learning from human feedback](#). *Preprint*,  
605 arXiv:2204.05862.

606 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,  
607 Amanda Askell, Jackson Kernion, Andy Jones, Anna  
608 Chen, Anna Goldie, Azalia Mirhoseini, Cameron

609 McKinnon, Carol Chen, Catherine Olsson, Christo-  
610 pher Olah, Danny Hernandez, Dawn Drain, Deep  
611 Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez,  
612 Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua  
613 Landau, Kamal Ndousse, Kamile Lukosuite, Liane  
614 Lovitt, Michael Sellitto, Nelson Elhage, Nicholas  
615 Schiefer, Noemi Mercado, Nova DasSarma, Robert  
616 Lasenby, Robin Larson, Sam Ringer, Scott John-  
617 ston, Shauna Kravec, Sheer El Showk, Stanislav Fort,  
618 Tamera Lanham, Timothy Telleen-Lawton, Tom Con-  
619 erly, Tom Henighan, Tristan Hume, Samuel R. Bow-  
620 man, Zac Hatfield-Dodds, Ben Mann, Dario Amodei,  
621 Nicholas Joseph, Sam McCandlish, Tom Brown, and  
622 Jared Kaplan. 2022b. [Constitutional ai: Harmless-  
623 ness from ai feedback](#). *Preprint*, arXiv:2212.08073.

624 Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec  
625 Koppel, Furong Huang, Dinesh Manocha, Am-  
626 rit Singh Bedi, and Mengdi Wang. 2024. [Maxmin-  
627 rlhf: Towards equitable alignment of large language  
628 models with diverse human preferences](#). *Preprint*,  
629 arXiv:2402.08925.

630 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa  
631 Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srini-  
632 vasan, Tianyi Zhou, Heng Huang, and Hongxia Jin.  
633 2024. [Alpagasus: Training a better alpaca model  
634 with fewer data](#). In *The Twelfth International Confer-  
635 ence on Learning Representations*.

636 Jonathan H Choi. 2024. [Measuring clarity in legal text](#).  
637 *U. Chi. L. Rev.*, 91:1.

638 Xu Chu, Ihab F. Ilyas, Sanjay Krishnan, and Jiannan  
639 Wang. 2016. [Data cleaning: Overview and emerging  
640 challenges](#). In *Proceedings of the 2016 International  
641 Conference on Management of Data, SIGMOD '16*,  
642 page 2201–2206, New York, NY, USA. Association  
643 for Computing Machinery.

644 Jeremy Cole, Michael Zhang, Daniel Gillick, Julian  
645 Eisenschlos, Bhuwan Dhingra, and Jacob Eisenstein.  
646 2023. [Selectively answering ambiguous questions](#).  
647 In *Proceedings of the 2023 Conference on Empiri-  
648 cal Methods in Natural Language Processing*, pages  
649 530–543, Singapore. Association for Computational  
650 Linguistics.

651 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and  
652 Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning  
653 of quantized llms](#). *Preprint*, arXiv:2305.14314.

654 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin,  
655 Shengding Hu, Zhiyuan Liu, Maosong Sun, and  
656 Bowen Zhou. 2023. [Enhancing chat language mod-  
657 els by scaling high-quality instructional conversa-  
658 tions](#). In *Proceedings of the 2023 Conference on  
659 Empirical Methods in Natural Language Processing*,  
660 pages 3029–3051, Singapore. Association for Com-  
661 putational Linguistics.

662 Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan  
663 Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng  
664 Zhang, KaShun SHUM, and Tong Zhang. 2023.

665	<a href="#">RAFT: Reward ranked finetuning for generative foundation model alignment.</a> <i>Transactions on Machine Learning Research.</i>	721
666		722
667		723
668	Romina Etezadi and Mehrnoush Shamsfard. 2023. The state of the art in open domain complex question answering: a survey. <i>Applied Intelligence</i> , 53(4):4124–4144.	724
669		725
670		726
671		727
672	H.A. Gleason. 1963. <i>Linguistics and English Grammar</i> . H.A. Gleason jr.	728
673		729
674	Benjamin M Gyori, Charles Tapley Hoyt, and Albert Steppi. 2022. Gilda: biomedical entity text normalization with machine-learned disambiguation as a service. <i>Bioinformatics Advances</i> , 2(1):vbac034.	730
675		731
676		732
677		733
678	Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023a. <a href="#">Beavertails: Towards improved safety alignment of LLM via a human-preference dataset.</a> In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	734
679		735
680		736
681		737
682		738
683		739
684		740
685	Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, Fanzhi Zeng, Kwan Yee Ng, Juntao Dai, Xuehai Pan, Aidan O’Gara, Yingshan Lei, Hua Xu, Brian Tse, Jie Fu, Stephen McAleer, Yaodong Yang, Yizhou Wang, Song-Chun Zhu, Yike Guo, and Wen Gao. 2023b. <a href="#">Ai alignment: A comprehensive survey.</a> <i>Preprint</i> , arXiv:2310.19852.	741
686		742
687		743
688		744
689		745
690		746
691		747
692		748
693		749
694	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. <a href="#">Mistral 7b.</a> <i>Preprint</i> , arXiv:2310.06825.	750
695		751
696		752
697		753
698		754
699		755
700		756
701		757
702	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. <a href="#">TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension.</a> In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.	758
703		759
704		760
705		761
706		762
707		763
708		764
709	Daniel Jurafsky. 1996. A probabilistic model of lexical and syntactic access and disambiguation. <i>Cognitive science</i> , 20(2):137–194.	765
710		766
711		767
712	Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joon-suk Park, and Jaewoo Kang. 2023a. <a href="#">Tree of clarifications: Answering ambiguous questions with retrieval-augmented large language models.</a> In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 996–1009, Singapore. Association for Computational Linguistics.	768
713		769
714		770
715		771
716		772
717		773
718		774
719	Najoung Kim, Phu Mon Htut, Samuel R. Bowman, and Jackson Petty. 2023b. <a href="#">(QA)<sup>2</sup>: Question answering with questionable assumptions.</a> In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8466–8487, Toronto, Canada. Association for Computational Linguistics.	775
720		776
		777
	Andreas K��pf, Yannic Kilcher, Dimitri von R��tte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Rich��rd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. <a href="#">Openassistant conversations - democratizing large language model alignment.</a> In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 47669–47681. Curran Associates, Inc.	
	Sushant Kumar, Sumit Datta, Vishakha Singh, Sanjay Kumar Singh, and Ritesh Sharma. 2024. Opportunities and challenges in data-centric ai. <i>IEEE Access</i> .	
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural questions: A benchmark for question answering research.</a> <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	
	Dongryeol Lee, Segwang Kim, Minwoo Lee, Hwanhee Lee, Joonsuk Park, Sang-Woo Lee, and Kyomin Jung. 2023. <a href="#">Asking clarification questions to handle ambiguity in open-domain QA.</a> In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11526–11544, Singapore. Association for Computational Linguistics.	
	Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. 2018. <a href="#">Scalable agent alignment via reward modeling: a research direction.</a> <i>CoRR</i> , abs/1811.07871.	
	Chin-Yew Lin and Franz Josef Och. 2004. <a href="#">Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics.</a> In <i>Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)</i> , pages 605–612, Barcelona, Spain.	
	Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha Swayamdipta, Noah Smith, and Yejin Choi. 2023. <a href="#">We’re afraid language models aren’t modeling ambiguity.</a> In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 790–807, Singapore. Association for Computational Linguistics.	
	Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2024a. <a href="#">What makes good data for alignment? a comprehensive study of automatic data se-</a>	

778		lection in instruction tuning. In <i>The Twelfth International Conference on Learning Representations</i> .	
779			
780	Zixuan Liu, Xiaolin Sun, and Zizhan Zheng. 2024b.	Enhancing llm safety via constrained direct preference optimization. <i>Preprint</i> , arXiv:2403.02475.	
781			
782			
783	Ilya Loshchilov and Frank Hutter. 2019.	Decoupled weight decay regularization. In <i>International Conference on Learning Representations</i> .	
784			
785			
786	Donald G Mackay and Thomas G Bever. 1967.	In search of ambiguity. <i>Perception &amp; Psychophysics</i> , 2:193–200.	
787			
788			
789	A. Majeed and S. Hwang. 2023.	Data-centric artificial intelligence, preprocessing, and the quest for transformative artificial intelligence systems development. <i>Computer</i> , 56(05):109–115.	
790			
791			
792			
793	Andrey Malinin and Mark Gales. 2021.	Uncertainty estimation in autoregressive structured prediction. In <i>International Conference on Learning Representations</i> .	
794			
795			
796			
797	Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022.	Peft: State-of-the-art parameter-efficient fine-tuning methods. <a href="https://github.com/huggingface/peft">https://github.com/huggingface/peft</a> .	
798			
799			
800			
801			
802	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020.	AmbigQA: Answering ambiguous open-domain questions. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5783–5797, Online. Association for Computational Linguistics.	
803			
804			
805			
806			
807			
808			
809	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022.	Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	
810			
811			
812			
813			
814			
815	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019.	Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 32.	
816			
817			
818			
819			
820			
821	Jonathan Pilault, Xavier Garcia, Arthur Bražinskas, and Orhan Firat. 2023.	Interactive-chain-prompting: Ambiguity resolution for crosslingual conditional generation with interaction. In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 455–483, Nusa Dua, Bali. Association for Computational Linguistics.	
822			
823			
824			
825			
826			
827			
828			
829			
830			
	Massimo Poesio and Ron Artstein. 2005.	The reliability of anaphoric annotation, reconsidered: Taking ambiguity into account. In <i>Proceedings of the Workshop on Frontiers in Corpus Annotations II: Pie in the Sky</i> , pages 76–83, Ann Arbor, Michigan. Association for Computational Linguistics.	831 832 833 834 835 836
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023.	Direct preference optimization: Your language model is secretly a reward model. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	837 838 839 840 841 842
	Sanford Schane. 2002.	Ambiguity and misunderstanding in the law. <i>T. Jefferson L. Rev.</i> , 25:167.	843 844
	Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022.	ASQA: Factoid questions meet long-form answers. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8273–8288, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	845 846 847 848 849 850 851
	Mark Stevenson and Yikun Guo. 2010.	Disambiguation in the biomedical domain: the role of ambiguity type. <i>Journal of biomedical informatics</i> , 43(6):972–981.	852 853 854
	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023.	Stanford alpaca: An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://github.com/tatsu-lab/stanford_alpaca</a> .	855 856 857 858 859
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023.	Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	860 861 862 863 864 865
	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2024.	Fine-tuning language models for factuality. In <i>The Twelfth International Conference on Learning Representations</i> .	866 867 868 869 870
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu,		871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888



889	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>Preprint</i> , arXiv:2307.09288.	945
890		946
891		947
892		948
893		949
894	Thomas Wasow, Amy Perfors, and David Beaver. 2005. The puzzle of ambiguity. <i>Morphology and the web of grammar: Essays in memory of Steven G. Lapointe</i> , pages 265–282.	950
895		951
896		952
897		953
898	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Trans-formers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	954
899		955
900		956
901		957
902		958
903		959
904		960
905		961
906		962
907		963
908		964
909		965
910	Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. <a href="#">Less: Selecting influential data for targeted instruction tuning</a> . <i>Preprint</i> , arXiv:2402.04333.	966
911		967
912		968
913		969
914	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. <a href="#">WizardLM: Empowering large pre-trained language models to follow complex instructions</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	970
915		971
916		972
917		973
918		974
919		975
920	Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. 2023. Alignment for honesty. <i>arXiv preprint arXiv:2312.07000</i> .	976
921		977
922		978
923	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. <a href="#">Do large language models know what they don't know?</a> In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8653–8665, Toronto, Canada. Association for Computational Linguistics.	979
924		980
925		981
926		982
927		983
928		984
929	Yuewei Yuan, Chaitanya Malaviya, and Mark Yatskar. 2023. <a href="#">AmbiCoref: Evaluating human and model sensitivity to ambiguous coreference</a> . In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 1023–1030, Dubrovnik, Croatia. Association for Computational Linguistics.	985
930		986
931		987
932		988
933		989
934		990
935	Lingxi Zhang, Jing Zhang, Xirui Ke, Haoyang Li, Xinmei Huang, Zhonghui Shao, Shulin Cao, and Xin Lv. 2023. A survey on complex factual question answering. <i>AI Open</i> , 4:1–12.	991
936		992
937		
938		
939	Michael Zhang and Eunsol Choi. 2021. <a href="#">SituatQA: Incorporating extra-linguistic contexts into QA</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7371–7387, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
940		
941		
942		
943		
944		
	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. <i>Advances in Neural Information Processing Systems</i> , 36.	
	<b>A Baseline Details</b>	
	In this section, we describe detailed implementation of the baselines.	
	<b>NAÏVE</b> We make a direct inference using template from Table 5. We evaluate the greedy generation result with temperature 0.	
	<b>AMBIGUITY-AWARE</b> We utilize the template from Table 6, where we explicitly describe how to handle ambiguity. Identically, we use the greedy generations for evaluation.	
	<b>SAMPLE REPETITION</b> Template from Table 5 is used to generate a single greedy generation and 10 sampled generations with sampling temperature 1.0. We measure the rate of sampled generations that match the greedy generation, which is reported to be best-calibrated (Cole et al., 2023). Samples with the measure below a specific threshold is considered ambiguous. We empirically select a threshold that demonstrates the best accuracy and F1-score with the least trade-off.	
	<b>SELF-ASK</b> We initially prompt the model with the template from Table 5 and generate a greedy generation. Then, the initial query and the generated answer is utilized with the template from Table 7 and prompt the model to verify the ambiguity of the query. We modified the prompt from Amayuelas et al. (2023) so that the model can specifically focus on ambiguity. The ambiguity detection is determined based on the model’s final verification.	
	<b>HONESTY-TUNED</b> The approach involves measuring the "expected accuracy" of sampled generations. The expected accuracy is measured by generating 10 samples with temperature 1.0 and measuring the average accuracy among the generations. Ground-truth ambiguous samples with expected accuracy below the specific threshold, in this case 0.1 from Yang et al. (2023), are labeled ambiguous. The samples classified as ambiguous is re-labeled with $y_{\text{clarify}}$ . The model $M$ is trained to generate the ground-truth $y$ given the input query $x$ and the inference template $t(\cdot)$ from Table 5.	
	<b>FULL-SET</b> The full training set is utilized for training. The ground-truth ambiguous samples are	



---

Answer the following question.  
 Question: <question>  
 Answer:

---

Table 5: Naïve prompting template.

---

Answer the following question. If the question is ambiguous, it is proper to answer with "The question is ambiguous".  
 Question: <question>  
 Answer:

---

Table 6: Ambiguity aware prompting. We explicitly describe how to handle ambiguity.

993 labeled with  $y_{\text{clarify}}$ . The model is trained to gener-  
 994 ate  $y$  given input question  $x$  and inference template  
 995  $t(\cdot)$  from Table 5.

996 **SUBSET** We randomly select  $|D|$  samples from  
 997 the training data, with the same number ( $|D|/2$ ) of  
 998 ambiguous and unambiguous samples. SUBSET is  
 999 trained in a same way as FULL-SET.

## 1000 B AmbigTriviaQA Construction Details

1001 AmbigTriviaQA is constructed by ambiguating the  
 1002 widely-used TriviaQA dataset (Joshi et al., 2017).  
 1003 We first prompt gpt-3.5-turbo to ambiguating the  
 1004 original question with the template from Table 8.  
 1005 To further validate the generation and control the  
 1006 quality of the dataset, we prompt gpt-3.5-turbo  
 1007 again for a secondary verification. We utilize the  
 1008 template in Table 9 and collect samples verified as  
 1009 ambiguous. This process yielded a total of 4,374  
 1010 question pairs to examine the model’s capability  
 1011 to interpret and generate responses to intentionally  
 1012 ambiguous queries. Examples from AmbigTrivi-  
 1013 aQA are demonstrated in Table 10.

## 1014 C Ambiguity Detection Details

1015 For ambiguous questions, we expect the model  
 1016 to generate clarification requests. Since there  
 1017 are various ways to express clarification requests,  
 1018 we use the following list of phrases to de-  
 1019 tect the requests. The presence of pre-defined  
 1020 ambiguity-related phrases in the model’s output is  
 1021 treated as a successful detection. The pre-defined  
 1022 phrases are the follows: [ambiguous, ambig,  
 1023 unclear, not clear, not sure, confused,  
 1024 confusing, vague, uncertain, doubtful,

---

Answer the following question. Given the question and answer, is the question ambiguous or unambiguous? Answer only ambiguous or unambiguous.  
 Question: <question>  
 Answer: <generated answer>

---

Is the question ambiguous or unambiguous? Answer only ambiguous or unambiguous.  
 Ambiguous or Unambiguous:

---

Table 7: Verification template for SELF-ASK. With the generated answer and the original question, the model is prompted to verify the ambiguity of the initial query.

---

Please make the following question ambiguous. Your task is to introduce ambiguity by altering the specificity of the noun phrase or omitting crucial details from the statement. Keep the rest of the sentence unchanged except for the modified sections. Generate only the revised statement.

Question: <question>  
 Ambiguation:

---

Table 8: Template to ambiguating the input query from TriviaQA. We prompt gpt-3.5-turbo for the generation.

doubt, questionable, clarify, not clear] 1025

## 1026 D Implementations Details

### 1027 D.1 Pipeline Details

1028 For explicit prediction (Stage 1), we utilize the  
 1029 same inference template as NAÏVE (Table 5) and  
 1030 the disambiguation is generated with the template  
 1031 from Table 11. We use the greedy generation for  
 1032 the disambiguated output. The threshold  $\epsilon$  is set to  
 1033 0.1 for filtering ambiguous inputs. For balancing  
 1034 training set size, if  $n > m$ , we randomly select  
 1035  $m$  samples from  $D_{\text{correct}}$ , where  $n = |D_{\text{correct}}|$  and  
 1036  $m = |D_{\text{ambig}}|$ . If  $n < m$ , we select  $n$  samples  
 1037 from  $D_{\text{ambig}}$  with the largest INFOGAIN. Finally,  
 1038 for  $y_{\text{clarify}}$ , we randomly select from the follow-  
 1039 ing pre-defined phrases : [The questions is  
 1040 ambiguous. Please clarify your question.  
 1041 Your question is ambiguous. Can you  
 1042 clarify your question? Your question is

---

An ambiguous question has multiple valid answers. Is the following question ambiguous with multiple possible answers? Answer only in Yes or No.

Question: <ambiguous generation>

Yes or No:

---

Table 9: Template for validating the generated ambiguous queries. We prompt gpt-3.5-turbo for the validation. Samples with the output "Yes" are considered a valid ambiguity and are selected as AmbigTriviaQA.

1043 not clear. Can you clarify your question  
1044 please?]

## 1045 **D.2 Training Details**

1046 For training, we applied AdamW optimizer  
1047 (Loshchilov and Hutter, 2019) with a batch size  
1048 of 32. We selected the model with the best perfor-  
1049 mance from learning rates {1e-3, 5e-4, 1e-4}  
1050 and training epochs {1, 2, 3}. All the exper-  
1051 iments were implemented with Pytorch (Paszke  
1052 et al., 2019) and Huggingface Transformers li-  
1053 brary (Wolf et al., 2020). For efficient training, we  
1054 applied QLoRA from Huggingface PEFT library  
1055 (Mangrulkar et al., 2022) with  $r=4$  and  $alpha=16$ .  
1056 The training takes about half an hour on a single  
1057 Tesla V100 GPU.

Original Question	Ambiguated Question
Which <b>volcano</b> in Tanzania is the highest mountain in Africa?	Which <b>geological formation</b> in Tanzania holds the title for the tallest landform in Africa?
What was <b>President Gerald Ford's</b> middle name?	What was the middle name of a <b>former U.S. president</b> ?
Where in England was <b>actor Nigel Hawthorne</b> born?	Where in the UK was <b>the actor</b> born?

Table 10: Example of the original question and its ambiguation from AmbigTriviaQA. The **ambiguated phrase** is highlighted in bold.

---

The question seems ambiguous, potentially requiring further details for clarity. Please disambiguate by providing specific context or constraints related to the question. This could include specifying any relevant time periods, locations, or other criteria necessary to narrow down possible interpretations. No additional specification is necessary for an unambiguous question.

Input Question: When did the frozen ride open at epcot?

Disambiguation: When did the frozen ride open at epcot?

Input Question: What is the legal age of marriage in usa?

Disambiguation: What is the legal age of marriage, without parental consent or other authorization, in all but two states in the usa?

Input Question: <question>

Disambiguation:

---

Table 11: Disambiguation template for implicit ambiguity measure. We provide 2-shot demonstration from AmbigQA train set.