

# Purging the Gray Zone: Latent-Geometric Denoising for Precise Knowledge Boundary Awareness

Anonymous ACL submission

## Abstract

Large language models (LLMs) often exhibit hallucinations due to their inability to accurately perceive their own knowledge boundaries. Existing abstention fine-tuning methods typically partition datasets directly based on response accuracy, causing models to suffer from severe label noise near the decision boundaries and consequently exhibit high rates of abstentions or hallucinations. This paper adopts a latent space representation perspective, revealing a “gray zone” near the decision hyperplane where internal belief ambiguity constitutes the core performance bottleneck. Based on this insight, we propose GEODE (Geometric Denoising) framework for abstention fine-tuning. This method constructs a truth hyperplane using linear probes and performs “geometric denoising” by employing geometric distance as a continuous abstention decision confidence metric. This approach filters out ambiguous boundary samples while retaining high-fidelity signals for fine-tuning. Experiments across multiple models (Llama3, Qwen3) and benchmark datasets (TriviaQA, NQ, SciQ, SimpleQA) demonstrate that GEODE significantly enhances model truthfulness and exhibits outstanding generalization capabilities in out-of-distribution (OOD) scenarios.

## 1 Introduction

Large Language Models (LLMs) have demonstrated outstanding performance across various natural language processing tasks (Grattafiori et al., 2024; Yang et al., 2025). However, it is universally acknowledged that LLMs exhibit hallucination—that is, generating responses that are factually inaccurate or fabricating answers (Zhang et al., 2025a). This issue underscores the urgent need to develop effective hallucination detection and mitigation methods.

A practicable solution to mitigate hallucination phenomena involves fine-tuning LLMs to remain

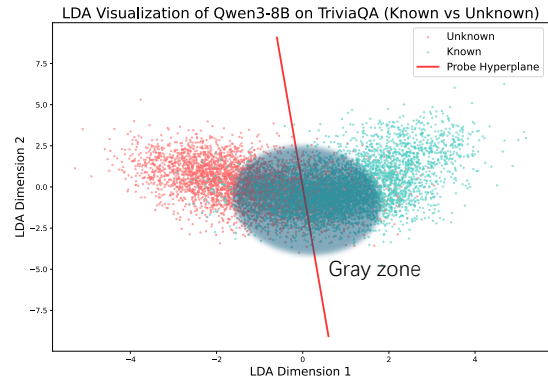


Figure 1: Visualization of the hidden states of questions that are known and unknown to the model. “Gray zone” refers to the overlapping area.

active when answering known questions and refrain from responding to queries beyond their knowledge scope (Wen et al., 2025; Li et al., 2025a). Specifically, these methods typically classify training data into “known” and “unknown” questions based on the correctness or accuracy of model responses, training models to answer the known set while replying “I don’t know” to the unknown set (Zhang et al., 2024a; Cheng et al., 2024). To reduce reliance on ground truth labels, alternative approaches utilize uncertainty metrics like semantic entropy to partition training data into known and unknown questions (Tjandra et al., 2024; Xue et al., 2025).

This type of method has demonstrated significant effectiveness across various models. However, current abstention fine-tuning often fails when internal confidence misaligns with external correctness and vague internal beliefs. Relying on stochastic accuracy to partition “known” and “unknown” sets introduces significant label noise—such as “lucky guesses” or formatting-driven failures. Training on these noisy heuristics forces the model to learn a jagged, contradictory decision boundary, ultimately leading to over-refusal or persistent hallucinations.

We employ *probing* method to analyze such cases. As shown in Figure 1, we visualize the known and unknown sets, with the red line representing the boundary obtained through training. It can be observed that the two sets exhibit significant overlap. Notably, while the central region exhibits overlap and boundary crossing, this issue does not occur in data points farther from the hyperplane. This indicates that a significant portion of the current abstention fine-tuning data contains noise. Therefore, a potential improvement is to discard noisy data and retain clean data in the training set so that the LLM can learn to distinguish the known from unknown cases more effectively.

To this end, we propose **Geometric Denoising** (GEODE) method for abstention fine-tuning, inspired by the linear representation hypothesis. This approach selects effective samples for abstention fine-tuning through linear probes, thereby enhancing LLMs’ self-awareness of their knowledge boundaries. The basic principle is, if the hidden state representation of a question is highly distant from the probe’s hyperplane, then the model can readily decide to reject or answer it. However, if it is very close to the hyperplane, then it indicates a difficult case to decide whether to reject or answer. Guided by this principle, we select samples with high probe confidence (distant from the probe hyperplane) and discard those with low probe confidence (close to the probe hyperplane).

Our main contributions are as follows:

- 1. Internal Representation Perspective:** We offer a novel diagnostic perspective on abstention fine-tuning by analyzing the latent space of LLMs. Our analysis reveals that suboptimal performance frequently stems from a “grey zone” near the latent decision boundary, where ambiguous representations introduce significant label noise.
- 2. Latent-Guided Denoised Dataset Curation:** We propose GEODE, a framework that leverages internal probes to curate high-quality fine-tuning datasets. By using geometric distance from the truthfulness hyperplane as a confidence metric, GEODE purges ambiguous boundary samples. This geometric denoising ensures the model trains on linearly separable signals, leading to sharper knowledge boundaries.
- 3. Empirical Superiority:** Extensive experiments across multiple architectures and benchmarks

demonstrate that GEODE significantly outperforms baselines. Our method also shows superior generalization in out-of-distribution (OOD) and abstention tasks.

## 2 Related Work

### 2.1 Hidden States of LLMs

Recent work suggests there is a “truthfulness” direction in latent space (Marks and Tegmark, 2024; Azaria and Mitchell, 2023). Liu et al. (2024) suggest there is a universal truthfulness hyperplane that exists within LLMs generalizing on cross-task, cross-domain, and in-domain. Some work probes the last token of a question to predict whether the model can answer it correctly without generating any tokens (Slobodkin et al., 2023; Snyder et al., 2024; Gottesman and Geva, 2024). To more effectively distinguish facts from errors, some work designs more complex features to train truthfulness probes and utilize information from model-generated answers (Orgad et al., 2025; Li et al., 2025b; Zhang et al., 2025b). Li et al. (2025b); Orgad et al. (2025). Truthfulness vectors can also be employed for hallucination mitigation by steering (Ji et al., 2025; Zhang et al., 2024b). Recent works suggest that models’ own internal judgments often lead to better overall factuality (Newman et al., 2025; Liang et al., 2024). In this work, we employ the truthfulness hyperplane as an internal confidence classifier to guide abstention fine-tuning for the model.

### 2.2 Abstention Fine-tuning

Abstention fine-tuning is a technique that teaches the model to abstain from answering questions that it does not know, while maintaining accuracy on known questions (Wen et al., 2025). Zhang et al. (2024a); Tjandra et al. (2024); Cheng et al. (2024) construct abstention-aware dataset based on whether the model can answer correctly, defining this as the model’s knowledge boundary. Then they fine-tune the model to refuse answering questions beyond its knowledge boundary while responding to those within it. Xu et al. (2024); Cheng et al. (2024); Brahman et al. (2024) use Direct Preference Optimization (DPO) (Rafailov et al., 2023) to train models to admit uncertainty when encountering unknown questions rather than outputting incorrect answers. Li et al. (2025c) employ adaptive contrastive learning to optimize LLMs’ abstention preferences. Huang et al. (2025); Cohen et al.

(2024) incorporate a dedicated “rejection” token into the model’s vocabulary and formulate an objective function that redistributes probability mass toward this token when the model is uncertain. Zheng et al. (2025); An and Xu (2025) train models to output binary confidence labels (“sure” vs. “unsure”) after generating an answer as a proxy for abstention, enabling them to reject low-confidence answers. Abstention fine-tuning may result in models being overly conservative or overly aggressive (Cheng et al., 2024; Zhu et al., 2025). In this work, we construct fine-tuning datasets based on the model’s internal beliefs to mitigate issues of over-rejection and over-hallucination.

### 3 Method

Our work aims to develop a hidden state approach to enhance the model’s awareness of its own knowledge boundaries. The main steps in our approach are to sample and arrange these diagnostic data into known/unknown subsets, train a hidden-state probe to quantify the model’s confidence in its responses, and finally perform targeted abstention fine-tuning on the curated samples to teach the model to abstain from answering questions outside its knowledge scope.

#### 3.1 Identify the Knowledge Boundary

First, the base LLM is prompted to answer every question within a source training dataset ( $D_0$ ). This acts as a diagnostic phase to examine what knowledge the model has actually internalized during its initial training. The model’s responses are then split into two distinct subsets based on the accuracy in answering each question.

**Known Knowledge ( $D_{0_{ik}}$ ):** These are the samples that the LLM correctly answered. The “ik” stands for “I know”. These question-answer pairs are retained in their original form to reinforce the retention of correct information during fine-tuning.

**Unknown Knowledge ( $D_{0_{idk}}$ ):** These are the samples that the LLM provided an incorrect response. For these cases, the original (incorrect) answer is discarded and replaced with a refusal message—specifically “I don’t know” in abstention fine-tuning (hence, “idk” for short).

#### 3.2 Curating Denoised Dataset

We train the probe model using  $D_{0_{ik}}$  and  $D_{0_{idk}}$  from the previous step. Specifically, we use the hidden states representation for the question-related

string as the feature  $\mathbf{x} = f_{\text{LLM}}(q)$ , and the correctness of the statement as a binary label  $y = \mathbb{I}(q \Rightarrow a) \in \{0, 1\}$ , to train a logistic regression probe, whose formula is:

$$f_{\text{probe}}(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b), \quad (1)$$

where  $\sigma$  is the sigmoid function, and  $w$  denotes the linear weight and  $b$  is the bias term. Next, we define the confidence of LLM answering question  $q$  with  $a$  by measuring the distance from  $\mathbf{x}$  to the learned hyperplane ( $\mathbf{w}; b$ ):

$$d(\mathbf{x}) = \frac{|\mathbf{w}^\top \mathbf{x} + b|}{\|\mathbf{w}\|_2} \quad (2)$$

$d(\mathbf{x})$  has a clear probabilistic meaning that reflects the model’s confidence in answering the given question.  $d(x) > 0$  means the model believes it can answer correctly, while  $d(x) < 0$  means the opposite, and a larger  $|d(x)|$  value indicates stronger confidence. It follows naturally that we can reformulate the goal of abstention fine-tuning as teaching the model to reject questions for which  $d(x) < 0$  and to answer those for which  $d(x) > 0$ . Accordingly, we partition the data into subsets of varying difficulty based on the magnitude of  $|d(x)|$ . Instead of using all available data, we select only samples farthest from the decision boundary, i.e., those with large  $|d(x)|$ , retaining the top  $X\%$  of samples for the fine-tuning task.

The detailed methods for extracting hidden state from LLMs,  $\mathbf{x} = f_{\text{LLM}}(q)$ , are as follows:

**Hidden state of the question (TBG):** We directly feed the question string  $q$  into the model and extract the hidden state of the last token of the question from the final layer, e.g., **token before generation** (TBG), which corresponds to the token immediately preceding the generation trigger for question encoding.

**Hidden state of the answer (SLT):** We first generate the model’s answer  $a$  via few-shot learning and greedy decoding (with temperature set to 0), and then feed concatenated sequence  $q \oplus a$  to the model, and retrieve the hidden state of the last token from the model’s final layer, e.g., **second last token** (SLT); this token captures the contextual representation of the entire question-answer sequence right before the end-of-sequence token.

#### 3.3 Abstention Fine-tuning on Subset

Given a dataset  $D^{\text{selected}}$  that consists of a known question set  $D_{\text{ik}}^{\text{selected}}$  and an unknown question set

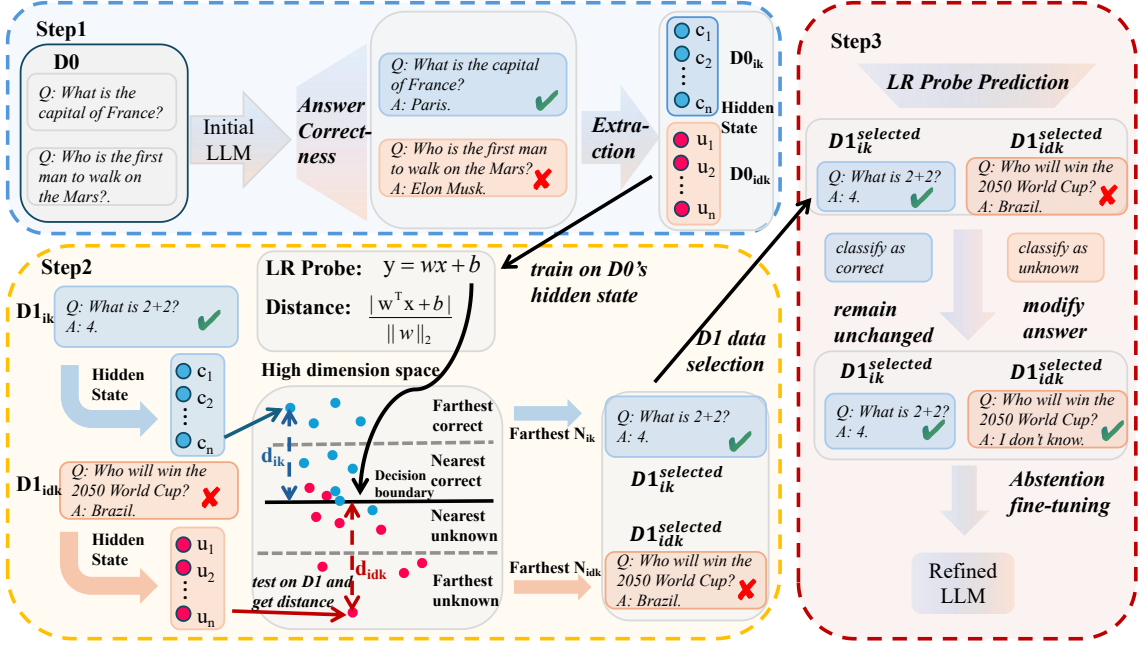


Figure 2: Overview of our method. GEODE contains three steps: (1) Identify the knowledge boundary by dividing the source data into two subsets ( $D0_{ik}$  and  $D0_{idk}$ ), and train a probe using these two subsets. (2) Calculate the distance from test data  $D1$  to the hyperplane learned by the probe, and select the subset of samples that are farthest from the hyperplane  $D1_{ik}^{selected}$  and  $D1_{idk}^{selected}$ . (3) Adjust the target answers based on the probe’s prediction results (retain correct answers for  $D1_{ik}^{selected}$  and replace incorrect answers with “I don’t know” for  $D1_{idk}^{selected}$ ), then conduct abstention fine-tuning on the selected subset.

$D_{idk}^{selected}$ , we modify the ground truth of  $D_{idk}^{selected}$  as “I don’t know” and keep the ground truth of  $D_{ik}^{selected}$  as the correct answer. Then we employ supervised fine-tuning with the cross-entropy loss:

$$\mathcal{L}_{(p_\theta)} = - \sum_{q \in D^{selected}} \sum_{t=1}^{|y^{(q)}|} \log p_\theta(y_t^{(q)} | \mathbf{I}, \mathbf{q}, \mathbf{y}_{t-1}^{(q)}), \quad (3)$$

in which  $p_\theta(y_t^{(q)})$  is the model’s predicted next-token probability distribution given the instruction ( $\mathbf{I}$ ), question ( $\mathbf{q}$ ), and the first  $t - 1$  tokens of the ground truth  $\mathbf{y}_{t-1}^{(q)}$ .

## 4 Experiments

### 4.1 Experimental Setting

**Datasets** We assess the effectiveness of GEODE on four open-ended question-answering tasks: **TriviaQA** (Joshi et al., 2017) which contains general knowledge QA pairs; **Natural Questions** (NQ) (Kwiatkowski et al., 2019) which contains questions from users’ queries to search engines; **SciQ** (Welbl et al., 2017) which contains science exam questions across multiple disciplines; and **SimpleQA** (Wei et al., 2024) which is adversarially collected against GPT-4 responses, using short-fact

---

### Algorithm 1 GEODE TBG Process

---

**Require:** Initial model  $M$ , QA dataset  $D_{src} = D0 + D1$ , percentage threshold  $X$

**Ensure:** Fine-tuned model  $M_{fine-tuned}$

- 1: **Step 1: Identify Knowledge Boundary**
  - 2: Test  $M$  on  $D0$
  - 3: Split  $D0$  into  $D0_{ik}$  and  $D0_{idk}$  by accuracy
  - 4: **Step 2: Curate Subsets Based on Distance**
  - 5: Get representation  $x = f_{LLM}(q)$ ,  $q \in D0$
  - 6: Train linear probe on  $D0$ :  $f_{probe}(x) = \sigma(w^\top x + b)$
  - 7: Compute distance on  $D1$ :  $d(x) = \frac{|w^\top x + b|}{\|w\|_2}$
  - 8: Determine threshold  $\theta$  as the  $X\%$ -th quantile of sorted  $d(x)$
  - 9: Select top  $X\%$  samples:  $D1^{selected} \leftarrow \{x \mid |d(x)| > \theta\}$
  - 10: **Step 3: Abstention Fine-tuning**
  - 11: Split  $D1^{selected}$  into  $D1_{ik}^{selected}$  ( $d(x) > 0$ ) and  $D1_{idk}^{selected}$  ( $d(x) < 0$ ).
  - 12: Replace ground truth of  $D1_{idk}^{selected}$  with “I don’t know.”
  - 13: Fine-tune  $M$  with cross-entropy loss on  $D1^{selected}$
  - 14: Return  $M_{fine-tuned}$
-

question-answering with single, indisputable answers to test whether the model truly “knows what it knows”. We use 10K samples from the TriviaQA training set for probe training, with the rest used for SFT. The validation split of TriviaQA is used for the in-domain (ID) test. We use NQ, SciQ, and SimpleQA for out-of-distribution (OOD) tests. Detailed evaluation dataset information is provided in Appendix A.

**Baselines** We compare GEODE with the following existing methods for abstention fine-tuning.

1. **IDK (Cheng et al., 2024)** directly prompts the model to abstain from uncertain questions.
2. **Uncertainty (Xu et al., 2025)** first prompts the model to answer questions as accurately as possible, then prompts the model to output binary uncertainty (sure or unsure).
3. **R-Tuning (Zhang et al., 2024a)** randomly selects a set of questions, categorizes them as known or unknown based on the model’s accuracy on each question, changes the ground truth for unknown questions to “I don’t know,” and then performs supervised fine-tuning.
4. **Probe-Tuning (Newman et al., 2025)**. The workflow of Probe-Tuning is identical to that of R-Tuning, with the key difference being that Probe-Pred utilizes the prediction results from truthful probes as the criteria for classifying questions into known and unknown categories.

GT \ R	Correctly answered	Wrongly answered	Abstained
Known	$N_1$	$N_2$	$N_3$
Unknown	–	$N_4$	$N_5$

Table 1: Abstention confusion matrix. R denotes the answer of the refined (fine-tuned) model. GT denotes the initial model. “Known” denotes that the initial model answered correctly, while “unknown” indicates that the initial model answered wrongly.

**Evaluation** For each test question, we classify the response as correct, wrong, or abstention. We use Llama3.1-8b-instruct (Grattafiori et al., 2024) as the judge to evaluate the correctness of answers generated by LLMs with 6-shot prompting. A response containing “I don’t know” is counted as an abstention. To measure the performance of abstention fine-tuning methods, we adopt three metrics, each reflecting unique aspects of performance,

based on the widely used abstention confusion matrix in Table 1. We use in-context learning (ICL) as the initial model.

**Helpfulness ( $F1_{\text{ans}}$ ):** For known questions, we calculate  $F1_{\text{ans}}$  (Kim et al., 2024) as the harmonic mean of answerable recall ( $\frac{N_1}{N_1+N_2+N_3}$ ) and answerable precision ( $\frac{N_1}{N_1+N_2+N_4}$ ).

**Truthfulness ( $F1_{\text{abs}}$ ):** For unknown questions, we calculate  $F1_{\text{abs}}$  (Kim et al., 2024) as the harmonic mean of unanswerable recall ( $\frac{N_5}{N_4+N_5}$ ) and unanswerable precision ( $\frac{N_5}{N_3+N_5}$ ).

**Reliability ( $F1_{\text{rel}}$ ):** Existing studies indicate that enhancing helpfulness leads to a decline in factuality (Xu et al., 2024; An and Xu, 2025). Therefore, we calculate the harmonic mean of metrics  $F1_{\text{ans}}$  and  $F1_{\text{abs}}$  as a reliability metric for comprehensive evaluation (An and Xu, 2025).

**Implementation Details** In this work, we choose Llama3-8B-Instruct (Grattafiori et al., 2024) and Qwen3-8B (Yang et al., 2025) as the initial models. We conduct experiments using SFT. We set the  $X$  as 25%. We employ logistic regression probe with L2 regularization. We use the Swift<sup>1</sup> framework to conduct fine-tuning using the AdamW optimizer, setting epoch to 3, learning rate to 1e-5, and batch size to 16. We use grid search to select the optimal hyperparameters. All the experiments are implemented on 4 Nvidia L40-48GB GPUs. During inference, we utilize the vLLM framework<sup>2</sup> to accelerate the process and employ a greedy search strategy to generate responses. Results are averaged over three different random seeds.

## 4.2 Main Results

We show the main experimental results of GEODE and all baseline methods in Table 2. The key observations from our experiments are:

**Effectiveness of GEODE** The fundamental limitation of existing baselines lies in their susceptibility to the “gray zone”—the latent region where internal belief is ambiguous and misaligned with external correctness. As illustrated in Figure 1, partitioning data based solely on response accuracy results in significant overlap between known and unknown representations. By utilizing the geometric distance from the truthfulness hyperplane, GEODE effectively purges this noise. Our results demonstrate that this denoising process allows the

<sup>1</sup><https://github.com/modelscope/ms-swift>

<sup>2</sup><https://github.com/vllm-project/vllm>

Dataset	TriviaQA			NQ			SciQ			SimpleQA		
Method	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>
<b>Llama3-8B-Instruct</b>												
<b>IDK</b>	78.7	35.3	48.8	<b>61.4</b>	56.3	58.7	<b>83.3</b>	5.6	10.4	14.6	64.0	23.8
<b>Uncertainty</b>	67.7	46.4	55.0	45.9	18.0	25.9	64.7	17.9	28.0	10.0	52.2	16.8
<b>R-Tuning</b>	77.3	71.6	74.4	47.0	78.1	58.7	69.9	58.2	63.5	14.6	<b>96.1</b>	25.4
<b>Probe-Tuning TBG</b>	78.8	72.8	75.7	51.6	<u>79.1</u>	62.5	81.4	53.4	64.5	12.8	<b>96.1</b>	22.6
<b>GEODE TBG</b>	<b>80.9</b>	<b>73.7</b>	<b>77.1</b>	<u>54.0</u>	<u>78.4</u>	<b>64.0</b>	81.9	58.5	<u>68.3</u>	<u>18.4</u>	94.8	30.7
<b>Probe-Tuning SLT</b>	75.8	71.3	73.4	44.8	78.3	56.9	75.4	58.9	66.1	<u>18.4</u>	95.8	30.9
<b>GEODE SLT</b>	<u>79.8</u>	<b>73.7</b>	<u>76.7</u>	52.6	<b>79.4</b>	<u>63.3</u>	<u>82.4</u>	<b>59.8</b>	<b>69.3</b>	<b>18.6</b>	96.0	<b>31.2</b>
<b>Qwen3-8B</b>												
<b>IDK</b>	74.0	55.1	63.2	<u>55.1</u>	65.4	59.8	81.8	44.5	57.6	11.7	58.6	19.5
<b>Uncertainty</b>	<u>75.9</u>	38.8	51.3	<b>57.3</b>	52.1	54.6	70.3	38.6	49.8	12.2	66.6	20.6
<b>R-Tuning</b>	75.8	74.3	75.0	50.3	70.5	58.7	<b>86.5</b>	45.1	<u>59.3</u>	12.6	63.6	21.0
<b>Probe-Tuning TBG</b>	75.0	71.7	73.3	51.2	70.5	59.4	<b>86.5</b>	44.2	58.5	12.6	63.4	21.0
<b>GEODE TBG</b>	<b>77.7</b>	<b>77.4</b>	<b>77.6</b>	54.3	<u>71.3</u>	<b>61.6</b>	81.8	<u>47.2</u>	59.8	<u>13.0</u>	<b>67.3</b>	<u>21.7</u>
<b>Probe-Tuning SLT</b>	75.4	72.8	74.1	50.1	69.2	58.1	85.7	45.8	<u>60.0</u>	12.3	62.0	20.6
<b>GEODE SLT</b>	75.6	<u>75.6</u>	<u>75.6</u>	51.6	<b>71.7</b>	<u>60.0</u>	84.6	<b>48.5</b>	<b>61.6</b>	<b>13.8</b>	65.9	<b>22.8</b>

Table 2: Performance on in-domain and out-of-domain question answering benchmarks. All results are multiplied by 100. The best result is bolded. The second best is underlined.

model to fine-tune on pure representations of its knowledge boundary. Consequently, GEODE consistently outperforms all baselines in the reliability metric ( $F1_{rel}$ ) across both Llama3-8B-Instruct and Qwen3-8B. For instance, on the TriviaQA in-domain task, GEODE (TBG) achieves a top  $F1_{rel}$  of 77.1 for Llama3 and 77.6 for Qwen, representing a significant margin over traditional R-Tuning and Probe-Tuning. GEODE also exhibits exceptional OOD generalization. On NQ, SciQ and SimpleQA, GEODE consistently maintains high  $F1_{rel}$  scores.

**Latent Representations: TBG vs. SLT** TBG and SLT exhibit comparable performance across benchmarks, with no single strategy consistently outperforming the other. While their absolute gains are similar, they reflect cognitive states at different stages: TBG captures the model’s initial confidence prior to output, whereas SLT incorporates the semantic context of the generated response. This result underscores the robustness of the GEODE framework regarding representation positioning. It demonstrates that geometric denoising effectively identifies and reduces noise regardless of the specific extraction point, confirming the universal utility of latent-geometric distance.

### 4.3 Evaluation on RAG Setting

**Experimental Setting** We evaluate the performance of the abstention methods in the Retrieval Augmented Generation (RAG) scenario. We use the RAG-Bench (Fang et al., 2024) dataset, which includes two settings: (1) **Golden**: golden retrieval, which contains contexts that include correct answers. (2) **Golden & RRN**: golden retrieval with relevant retrieval noise, which includes golden retrieval and context relevant to the question statement but lacks the correct answers. We feed the context alongside the question into the model.

**Experimental Results** Table 3 shows the performance in the RAG setting. Similar to the main experimental results, IDK and R-Tuning performed worse than abstention fine-tuning based on probing. Overall, our method achieves the best performance relative to the baseline. Within the SLT-based setting, our method not only achieves superior accuracy but also yields a lower hallucination rate compared to Probe-Tuning. This dual improvement highlights the efficacy of geometric denoising in establishing more reliable knowledge boundaries. Our method maintains the highest reliability ( $F1_{rel}$ ) even under noisy retrieval scenarios, demonstrating its robustness in practical applications.

	Golden					Golden & RRN				
	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>	Acc.	Hallu.	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>	Acc.	Hallu.
ICL	-	-	-	71.2	28.8	-	-	-	63.8	36.2
IDK	76.1	45.5	57.0	58.1	17.9	71.6	45.6	55.7	52.7	23.6
R-Tuning	80.8	52.5	63.7	61.0	11.8	74.4	53.3	62.1	50.4	15.3
Probe-tuning TBG	83.8	47.9	61.0	67.3	13.0	81.1	48.6	60.8	63.3	16.8
GEODE TBG	75.2	52.5	61.8	51.5	9.2	74.3	56.9	64.5	49.1	11.4
Probe-Tuning SLT	79.1	52.7	63.2	54.8	11.9	75.9	54.7	63.5	50.0	15.7
GEODE SLT	79.9	61.1	69.3	55.8	9.2	75.8	57.9	65.7	53.4	13.0

Table 3: RAG results. Acc. is accuracy. Hallu. is hallucination rate.

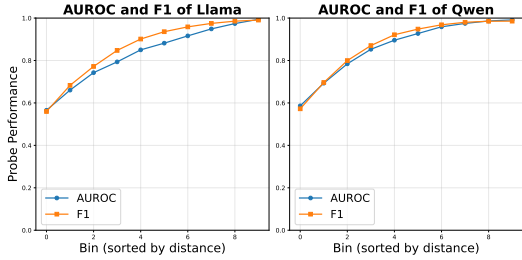


Figure 3: Probing performance vs. distance. Bin0 denotes the nearest subset, and bin9 the farthest subset. Accuracy and F1 score increase as distance increases.

## 5 Analysis

### 5.1 Ablation Studies

To validate our hypothesis that the geometric distance to the probing hyperplane serves as a reliable proxy for sample quality, we partitioned the training data into three tiers: Ours-Farthest (75%-100%), Ours-Middle (37.5%-62.5%), and Ours-Nearest (0-25%). As shown in Table 4, a clear performance gradient is observed across all metrics. Specifically, on Qwen3-8B (TriviaQA), Ours-Farthest achieves an F1<sub>rel</sub> of 77.6, markedly surpassing the Middle (74.6) and Nearest (73.2) tiers. This performance decay is inherently rooted in the noise density of the decision boundary: samples located near the hyperplane (the “grey zone”) represent instances where the model’s internal representations for known and unknown knowledge are highly entangled. Fine-tuning on these ambiguous samples introduces contradictory gradients that blur the model’s knowledge boundaries. By selectively training on the “farthest” samples, GEODE effectively performs geometric denoising, ensuring that the model learns from high-confidence, linearly separable signals. To further illustrate this distance-quality correlation, we bin the training

Dataset	TriviaQA			NQ		
Metric	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>	F1 <sub>ans</sub>	F1 <sub>abs</sub>	F1 <sub>rel</sub>
Ours TBG	<b>77.7</b>	<b>77.4</b>	<b>77.6</b>	<b>54.3</b>	71.3	<b>61.6</b>
-middle	74.9	74.2	74.6	50.0	70.2	58.4
-nearest	75.2	71.2	73.2	49.7	69.6	58.0
Ours SLT	75.6	75.6	75.6	51.6	71.7	60.0
-middle	75.4	72.8	74.1	50.1	69.2	58.1
-nearest	73.8	67.5	70.5	48.9	67.2	56.6
R-Tuning-01	75.5	77.0	76.2	51.6	<b>71.8</b>	60.0

Table 4: Ablation study on distance-based data partitioning for abstention fine-tuning. Samples are partitioned into three tiers based on their geometric distance  $|d(x)|$  to the probing hyperplane.

data into ten equal subsets sorted by distance in ascending order. As visualized in Figure 3, the probe’s predictive accuracy scales monotonically with distance. Notably, the AUROC for samples nearest to the hyperplane drops below 0.6, approaching random chance, confirming that proximity to the boundary is a primary source of aleatoric noise. In contrast, distal samples provide a clean signal for knowledge boundaries.

We further compare this geometric denoising approach with R-Tuning-01, which conceptually shares a similar objective by training only on samples with 100% response consistency (either all correct or all incorrect) across 10 times of independent samples. While R-Tuning-01 attempts to filter ambiguity via external output stability, it remains inferior to our method.

### 5.2 Effects of Positive Proportions in Training

We fix the training set size at 20,000 and conduct experiments under different positive sample pro-

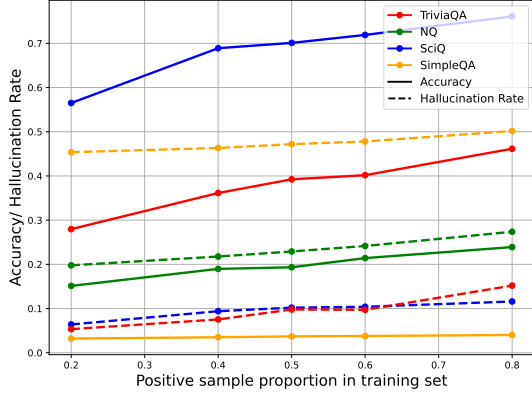


Figure 4: Accuracy and hallucination rate according to different positive proportions in the training set.

portions. As shown in Figure 4, both accuracy and hallucination rates increase as the proportion of positive samples in the training set go up. High-accuracy datasets like TriviaQA and SciQ exhibit significant sensitivity to changes in the positive-to-negative ratio within the training set, whereas NQ and SimpleQA show less pronounced variation. This suggests that excessively high negative sample ratios may cause models to over-abstain from known questions. The results also imply that abstention fine-tuning cannot eliminate over-abstention and hallucination simultaneously.

### 5.3 Evaluation on Unanswerable Datasets

**Experimental Setup** To evaluate the model’s ability for identifying unanswerable queries, we test performance on three specialized benchmarks: (1) **Alcuna** (Yin et al., 2023a), containing synthetic entity-based questions; (2) **FalseQA** (Hu et al., 2023), containing questions contradicting common sense; and (3) **Self-Aware (SA)** (Yin et al., 2023b), containing inherently unanswerable questions.

**Performance Analysis** As presented in Table 5, GEODE demonstrates highly competitive abstention rates, particularly with the TBG variant. While performance varies across benchmarks, our method shows notable strengths in specific scenarios. For instance, on the Alcuna dataset using Qwen3-8B, GEODE (TBG) achieves an abstention rate of 94.5%, a substantial improvement over the 78.2% reached by R-Tuning. We found that the rejection rate on Alcuna was higher than that on FalseQA and SA. Notably, overall rejection rates are higher on Alcuna than on FalseQA and SA. We attribute this to the nature of entities: while FalseQA and SA involve real-world concepts, Alcuna uses synthetic

Method	Alcuna	FalseQA	SA
<b>Llama3-8b-Instruction</b>			
<b>IDK</b>	78.4	49.8	50.9
<b>Uncertainty</b>	19.2	8.5	15.4
<b>R-Tuning</b>	98.3	91.8	98.1
<b>Probe-Tuning TBG</b>	99.2	91.8	<b>99.7</b>
<b>GEODE TBG</b>	98.2	92.4	99.2
<b>Probe-Tuning SLT</b>	<b>99.7</b>	<b>95.3</b>	<b>99.7</b>
<b>GEODE SLT</b>	99.0	90.6	98.7
<b>Qwen3-8B</b>			
<b>IDK</b>	85.9	62.4	74.4
<b>Uncertainty</b>	38.5	24.5	39.6
<b>R-Tuning</b>	78.2	67.5	82.6
<b>Probe-Tuning TBG</b>	84.0	<b>72.2</b>	84.1
<b>GEODE TBG</b>	<b>94.5</b>	67.6	<b>88.3</b>
<b>Probe-Tuning SLT</b>	69.7	65.2	80.6
<b>GEODE SLT</b>	94.3	70.5	86.9

Table 5: Abstention rates (%) across specialized benchmarks. High scores indicate effective identification of unanswerable or deceptive queries.

ones. This suggests that familiar domains trigger an illusion of knowledge, where the presence of known entities induces generative overconfidence. Consequently, models struggle more to recognize their epistemic limits in familiar contexts than in entirely novel, synthetic ones.

## 6 Conclusion

In this work, we introduce GEODE, a novel framework for abstention fine-tuning that leverages the latent geometry of LLMs. Moving beyond traditional methods that rely solely on external response accuracy, we propose a diagnostic perspective by analyzing the model’s internal representation space. We found a critical “grey zone” near the latent decision hyperplane, where ambiguous internal beliefs introduce significant noise that hinders a model’s ability to perceive its own knowledge boundaries. By employing geometric denoising, GEODE systematically purges these ambiguous boundary samples, ensuring that the model is fine-tuned on high-fidelity, linearly separable signals. Extensive experiments across multiple models (Llama3, Qwen3) and benchmarks demonstrate that our approach significantly enhances model reliability and exhibits superior generalization in OOD, RAG, and deceptive scenarios.

## 525 Limitations

526 Our method has been validated on models with up  
527 to 8B parameters; however, it has not been eval-  
528 uated on models with larger parameter sizes. We  
529 employ linear representations to model truthful-  
530 ness, but linear methods may not fully capture the  
531 complexity of truth-related behaviors, and a multi-  
532 dimensional framework is necessary for more ac-  
533 curate modeling in larger models as suggested by  
534 recent evidence (Yu et al., 2025). Additionally, our  
535 experiments have yet to include reasoning tasks  
536 or long-form generation, which are critical for fur-  
537 ther evaluating the robustness and scalability of  
538 our approach in addressing more intricate linguis-  
539 tic challenges. Although this work employs highly  
540 discriminative data for fine-tuning, abstention fine-  
541 tuning carries an unavoidable risk: over-abstention,  
542 which may cause the model to refuse to output  
543 answers it was highly confident about prior to fine-  
544 tuning. Future research may need to continuously  
545 teach models to express uncertainty during the pre-  
546 training stage.

## 547 References

548 Hao An and Yang Xu. 2025. [Teaching llms to ab-](#)  
549 [stain via fine-grained semantic confidence reward.](#)  
550 *Preprint*, arXiv:2510.24020.

551 Amos Azaria and Tom Mitchell. 2023. [The internal](#)  
552 [state of an LLM knows when it’s lying.](#) In *The 2023*  
553 *Conference on Empirical Methods in Natural Lan-*  
554 *guage Processing*.

555 Faeze Brahman, Sachin Kumar, Vidhisha Balachan-  
556 dran, Pradeep Dasigi, Valentina Pyatkin, Abhilasha  
557 Ravichander, Sarah Wiegrefe, Nouha Dziri, Khy-  
558 athi Chandu, Jack Hessel, Yulia Tsvetkov, Noah A.  
559 Smith, Yejin Choi, and Hannaneh Hajishirzi. 2024.  
560 [The art of saying no: Contextual noncompliance in](#)  
561 [language models.](#) In *The Thirty-eight Conference on*  
562 *Neural Information Processing Systems Datasets and*  
563 *Benchmarks Track*.

564 Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wen-  
565 wei Zhang, Zhangyue Yin, Shimin Li, Linyang Li,  
566 Zhengfu He, Kai Chen, and Xipeng Qiu. 2024. [Can](#)  
567 [AI assistants know what they don’t know?](#) In *Forty-*  
568 *first International Conference on Machine Learning*.

569 Roi Cohen, Konstantin Dobler, Eden Biran, and Gerard  
570 de Melo. 2024. [I don’t know: Explicit modeling](#)  
571 [of uncertainty with an \[idk\] token.](#) In *Advances in*  
572 *Neural Information Processing Systems*, volume 37,  
573 pages 10935–10958. Curran Associates, Inc.

574 Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiao-  
575 jun Chen, and Ruifeng Xu. 2024. [Enhancing noise](#)

[robustness of retrieval-augmented language models](#)  
[with adaptive adversarial training.](#) In *Proceedings*  
*of the 62nd Annual Meeting of the Association for*  
*Computational Linguistics (Volume 1: Long Papers)*,  
pages 10028–10039, Bangkok, Thailand. Association  
for Computational Linguistics.

Daniela Gottesman and Mor Geva. 2024. [Estimating](#)  
[knowledge in large language models without gen-](#)  
[erating a single token.](#) In *Proceedings of the 2024*  
*Conference on Empirical Methods in Natural Lan-*  
*guage Processing*, pages 3994–4019, Miami, Florida,  
USA. Association for Computational Linguistics.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,  
Abhinav Pandey, Abhishek Kadian, Ahmad Al-  
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-  
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh  
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-  
tra, Archie Sravankumar, Artem Korenev, Arthur  
Hinsvark, and 542 others. 2024. [The llama 3 herd of](#)  
[models.](#) *Preprint*, arXiv:2407.21783.

Shengding Hu, Yifan Luo, Huadong Wang, Xingyi  
Cheng, Zhiyuan Liu, and Maosong Sun. 2023. [Won’t](#)  
[get fooled again: Answering questions with false](#)  
[premises.](#) In *Proceedings of the 61st Annual Meet-*  
*ing of the Association for Computational Linguistics*  
*(Volume 1: Long Papers)*, pages 5626–5643, Toronto,  
Canada. Association for Computational Linguistics.

Lei Huang, Xiaocheng Feng, Weitao Ma, Yuchun Fan,  
Xiachong Feng, Yuxuan Gu, Yangfan Ye, Liang Zhao,  
Weihong Zhong, Baoxin Wang, Dayong Wu, Guop-  
ing Hu, Lingpeng Kong, Tong Xiao, Ting Liu, and  
Bing Qin. 2025. [Alleviating hallucinations from](#)  
[knowledge misalignment in large language models](#)  
[via selective abstention learning.](#) In *Proceedings*  
*of the 63rd Annual Meeting of the Association for*  
*Computational Linguistics (Volume 1: Long Papers)*,  
pages 24564–24579, Vienna, Austria. Association  
for Computational Linguistics.

Ziwei Ji, Lei Yu, Yeskendir Koishekenov, Yejin Bang,  
Anthony Hartshorn, Alan Schelten, Cheng Zhang,  
Pascale Fung, and Nicola Cancedda. 2025. [Calibrat-](#)  
[ing verbal uncertainty as a linear feature to reduce](#)  
[hallucinations.](#) In *Proceedings of the 2025 Confer-*  
*ence on Empirical Methods in Natural Language*  
*Processing*, pages 3769–3793, Suzhou, China. Asso-  
ciation for Computational Linguistics.

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke  
Zettlemoyer. 2017. [TriviaQA: A large scale distantly](#)  
[supervised challenge dataset for reading comprehen-](#)  
[sion.](#) In *Proceedings of the 55th Annual Meeting of*  
*the Association for Computational Linguistics (Vol-*  
*ume 1: Long Papers)*, pages 1601–1611, Vancouver,  
Canada. Association for Computational Linguistics.

Hyuhng Joon Kim, Youna Kim, Cheonbok Park, Jun-  
yeob Kim, Choonghyun Park, Kang Min Yoo, Sang-  
goo Lee, and Taek Kim. 2024. [Aligning language](#)  
[models to explicitly handle ambiguity.](#) In *Proceed-*  
*ings of the 2024 Conference on Empirical Methods*

576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633

634	<i>in Natural Language Processing</i> , pages 1989–2007,	<i>The Thirteenth International Conference on Learning</i>	691
635	Miami, Florida, USA. Association for Computational	<i>Representations</i> .	692
636			
637	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	693
638	field, Michael Collins, Ankur Parikh, Chris Alberti,	pher D Manning, Stefano Ermon, and Chelsea Finn.	694
639	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	2023. <a href="#">Direct preference optimization: Your language</a>	695
640	ton Lee, Kristina Toutanova, Llion Jones, Matthew	<a href="#">model is secretly a reward model</a> . In <i>Advances in</i>	696
641	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	<i>Neural Information Processing Systems</i> , volume 36,	697
642	Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natu-</a>	pages 53728–53741. Curran Associates, Inc.	698
643	<a href="#">ral questions: A benchmark for question answering</a>		
644	<a href="#">research</a> . <i>Transactions of the Association for Compu-</i>	Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido	699
645	<i>tational Linguistics</i> , 7:452–466.	Dagan, and Shauli Ravfogel. 2023. <a href="#">The curious case</a>	700
		<a href="#">of hallucinatory (un)answerability: Finding truths in</a>	701
646	Moxin Li, Yong Zhao, Wenxuan Zhang, Shuaiyi Li,	<a href="#">the hidden states of over-confident large language</a>	702
647	Wenya Xie, See-Kiong Ng, Tat-Seng Chua, and Yang	<a href="#">models</a> . In <i>Proceedings of the 2023 Conference on</i>	703
648	Deng. 2025a. <a href="#">Knowledge boundary of large lan-</a>	<i>Empirical Methods in Natural Language Processing</i> ,	704
649	<a href="#">guage models: A survey</a> . In <i>Proceedings of the</i>	pages 3607–3625, Singapore. Association for Com-	705
650	<i>63rd Annual Meeting of the Association for Compu-</i>	putational Linguistics.	706
651	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages		
652	5131–5157, Vienna, Austria. Association for Compu-	Ben Snyder, Marius Moisescu, and Muhammad Bilal	707
653	tational Linguistics.	Zafar. 2024. <a href="#">On early detection of hallucinations in</a>	708
		<a href="#">factual question answering</a> . In <i>Proceedings of the</i>	709
654	Qing Li, Jiahui Geng, Zongxiong Chen, Derui Zhu,	<i>30th ACM SIGKDD Conference on Knowledge Dis-</i>	710
655	Yuxia Wang, Congbo Ma, Chenyang Lyu, and Fakhri	<i>covery and Data Mining, KDD '24</i> , page 2721–2732,	711
656	Karray. 2025b. <a href="#">HD-NDEs: Neural differential equa-</a>	New York, NY, USA. Association for Computing	712
657	<a href="#">tions for hallucination detection in LLMs</a> . In <i>Pro-</i>	Machinery.	713
658	<i>ceedings of the 63rd Annual Meeting of the Associa-</i>		
659	<i>tion for Computational Linguistics (Volume 1: Long</i>	Benedict Aaron Tjandra, Muhammed Razzak, Jannik	714
660	<i>Papers)</i> , pages 6173–6186, Vienna, Austria. Associa-	Kossen, Kunal Handa, and Yarin Gal. 2024. <a href="#">Fine-</a>	715
661	tion for Computational Linguistics.	<a href="#">tuning large language models to appropriately abstain</a>	716
		<a href="#">with semantic entropy</a> . In <i>Neurips Safe Generative</i>	717
662	Yinghui Li, Haojing Huang, Jiayi Kuang, Yangning	<i>AI Workshop 2024</i> .	718
663	Li, Shu-Yu Guo, Chao Qu, Xiaoyu Tan, Hai-Tao		
664	Zheng, Ying Shen, and Philip S. Yu. 2025c. <a href="#">Refine</a>	Jason Wei, Nguyen Karina, Hyung Won Chung,	719
665	<a href="#">knowledge of large language models via adaptive</a>	Yunxin Joy Jiao, Spencer Papay, Amelia Glaese,	720
666	<a href="#">contrastive learning</a> . In <i>The Thirteenth International</i>	John Schulman, and William Fedus. 2024. <a href="#">Mea-</a>	721
667	<i>Conference on Learning Representations</i> .	<a href="#">suring short-form factuality in large language models</a> .	722
		<i>Preprint</i> , arXiv:2411.04368.	723
668	Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiax-		
669	ing Zhang. 2024. <a href="#">Learning to trust your feelings:</a>	Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017.	724
670	<a href="#">Leveraging self-awareness in llms for hallucination</a>	<a href="#">Crowdsourcing multiple choice science questions</a> .	725
671	<a href="#">mitigation</a> . <i>Preprint</i> , arXiv:2401.15449.	In <i>Proceedings of the 3rd Workshop on Noisy User-</i>	726
		<i>generated Text</i> , pages 94–106, Copenhagen, Den-	727
672	Junteng Liu, Shiqi Chen, Yu Cheng, and Junxian He.	mark. Association for Computational Linguistics.	728
673	2024. <a href="#">On the universal truthfulness hyperplane inside</a>		
674	<a href="#">LLMs</a> . In <i>Proceedings of the 2024 Conference on</i>	Bingbing Wen, Jihan Yao, Shangbin Feng, Chenjun Xu,	729
675	<i>Empirical Methods in Natural Language Processing</i> ,	Yulia Tsvetkov, Bill Howe, and Lucy Lu Wang. 2025.	730
676	pages 18199–18224, Miami, Florida, USA. Associa-	<a href="#">Know your limits: A survey of abstention in large</a>	731
677	tion for Computational Linguistics.	<a href="#">language models</a> . <i>Transactions of the Association for</i>	732
		<i>Computational Linguistics</i> , 13:529–556.	733
678	Samuel Marks and Max Tegmark. 2024. <a href="#">The geometry</a>		
679	<a href="#">of truth: Emergent linear structure in large language</a>	Chenjun Xu, Bingbing Wen, Bin Han, Robert Wolfe,	734
680	<a href="#">model representations of true/false datasets</a> . In <i>First</i>	Lucy Lu Wang, and Bill Howe. 2025. <a href="#">Do language</a>	735
681	<i>Conference on Language Modeling</i> .	<a href="#">models mirror human confidence? exploring psycho-</a>	736
		<a href="#">logical insights to address overconfidence in LLMs</a> .	737
682	Benjamin Newman, Abhilasha Ravichander, Jaehun	In <i>Findings of the Association for Computational</i>	738
683	Jung, Rui Xin, Hamish Ivison, Yegor Kuznetsov,	<i>Linguistics: ACL 2025</i> , pages 25655–25672, Vienna,	739
684	Pang Wei Koh, and Yejin Choi. 2025. <a href="#">The curious</a>	Austria. Association for Computational Linguistics.	740
685	<a href="#">case of factuality finetuning: Models' internal beliefs</a>		
686	<a href="#">can improve factuality</a> . <i>Preprint</i> , arXiv:2507.08371.	Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai	741
687	Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Re-	Fan, Lu Chen, and Kai Yu. 2024. <a href="#">Rejection improves</a>	742
688	ichart, Idan Szepkter, Hadas Kotek, and Yonatan Be-	<a href="#">reliability: Training LLMs to refuse unknown ques-</a>	743
689	linkov. 2025. <a href="#">LLMs know more than they show: On</a>	<a href="#">tions using RL from knowledge feedback</a> . In <i>First</i>	744
690	<a href="#">the intrinsic representation of LLM hallucinations</a> . In	<i>Conference on Language Modeling</i> .	745

746	Boyang Xue, Fei Mi, Qi Zhu, Hongru Wang, Rui Wang, Sheng Wang, Erxin Yu, Xuming Hu, and Kam-Fai Wong. 2025. <a href="#">UAlign: Leveraging uncertainty estimations for factuality alignment on large language models</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6002–6024, Vienna, Austria. Association for Computational Linguistics.	804																																				
747		805																																				
748		806																																				
749		807																																				
750		808																																				
751																																						
752		809																																				
753		810																																				
754	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. <a href="#">Qwen3 technical report</a> . <i>Preprint</i> , arXiv:2505.09388.	811																																				
755		812																																				
756		813																																				
757																																						
758		814																																				
759		815																																				
760		816																																				
761	Xunjian Yin, Baizhou Huang, and Xiaojun Wan. 2023a. <a href="#">ALCUNA: Large language models meet new knowledge</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1397–1414, Singapore. Association for Computational Linguistics.	817																																				
762		818																																				
763		819																																				
764		820																																				
765		821																																				
766																																						
767	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023b. <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.	822																																				
768		823																																				
769		824																																				
770		825																																				
771		826																																				
772																																						
773	Stanley Yu, Vaidehi Bulusu, Oscar Yasunaga, Clayton Lau, Cole Blondin, Sean O'Brien, Kevin Zhu, and Vasu Sharma. 2025. <a href="#">From directions to cones: Exploring multidimensional representations of propositional facts in llms</a> . <i>Preprint</i> , arXiv:2505.21800.																																					
774																																						
775																																						
776																																						
777																																						
778	Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024a. <a href="#">R-tuning: Instructing large language models to say 'I don't know'</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 7113–7139, Mexico City, Mexico. Association for Computational Linguistics.																																					
779																																						
780																																						
781																																						
782																																						
783																																						
784																																						
785																																						
786																																						
787	Shaolei Zhang, Tian Yu, and Yang Feng. 2024b. <a href="#">TruthX: Alleviating hallucinations by editing large language models in truthful space</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8908–8949, Bangkok, Thailand. Association for Computational Linguistics.																																					
788																																						
789																																						
790																																						
791																																						
792																																						
793																																						
794	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Chen Xu, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2025a. <a href="#">Siren's song in the ai ocean: A survey on hallucination in large language models</a> . <i>Preprint</i> , arXiv:2309.01219.																																					
795																																						
796																																						
797																																						
798																																						
799																																						
800																																						
801	Zhenliang Zhang, Xinyu Hu, Huixuan Zhang, Junzhe Zhang, and Xiaojun Wan. 2025b. <a href="#">ICR probe: Tracking hidden state dynamics for reliable hallucination detection in LLMs</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 17986–18002, Vienna, Austria. Association for Computational Linguistics.																																					
802																																						
803																																						
	Hang Zheng, Hongshen Xu, Yuncong Liu, Shuai Fan, Lu Chen, Pascale Fung, and Kai Yu. 2025. <a href="#">Enhancing LLM reliability via explicit knowledge boundary modeling</a> . In <i>Second Conference on Language Modeling</i> .																																					
	Runchuan Zhu, Xinke Jiang, Jiang Wu, Zhipeng Ma, Jiahe Song, Fengshuo Bai, Dahua Lin, Lijun Wu, and Conghui He. 2025. <a href="#">GRAIT: Gradient-driven refusal-aware instruction tuning for effective hallucination mitigation</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 4006–4021, Albuquerque, New Mexico. Association for Computational Linguistics.																																					
	<b>A Dataset Details</b>																																					
	The table displays statistical information about the datasets, where SimpleQA and NQ are more challenging datasets, while TriviaQA and NQ are relatively simpler. Alcuna, FalseQA and Self-Aware contain unanswerable questions.																																					
	<table border="1"> <thead> <tr> <th></th> <th>Size</th> <th>Llama Acc.</th> <th>Qwen Acc.</th> </tr> </thead> <tbody> <tr> <td><b>TriviaQA Train</b></td> <td>87622</td> <td>66.17</td> <td>54.93</td> </tr> <tr> <td><b>TriviaQA Val</b></td> <td>11313</td> <td>65.55</td> <td>54.65</td> </tr> <tr> <td><b>NQ</b></td> <td>3610</td> <td>40.86</td> <td>34.04</td> </tr> <tr> <td><b>SciQ</b></td> <td>1000</td> <td>72.50</td> <td>81.60</td> </tr> <tr> <td><b>SimpleQA</b></td> <td>4326</td> <td>6.80</td> <td>5.59</td> </tr> <tr> <td><b>Alcuna</b></td> <td>2001</td> <td>-</td> <td>-</td> </tr> <tr> <td><b>FalseQA</b></td> <td>1374</td> <td>-</td> <td>-</td> </tr> <tr> <td><b>Self-Aware</b></td> <td>1032</td> <td>-</td> <td>-</td> </tr> </tbody> </table>		Size	Llama Acc.	Qwen Acc.	<b>TriviaQA Train</b>	87622	66.17	54.93	<b>TriviaQA Val</b>	11313	65.55	54.65	<b>NQ</b>	3610	40.86	34.04	<b>SciQ</b>	1000	72.50	81.60	<b>SimpleQA</b>	4326	6.80	5.59	<b>Alcuna</b>	2001	-	-	<b>FalseQA</b>	1374	-	-	<b>Self-Aware</b>	1032	-	-	
	Size	Llama Acc.	Qwen Acc.																																			
<b>TriviaQA Train</b>	87622	66.17	54.93																																			
<b>TriviaQA Val</b>	11313	65.55	54.65																																			
<b>NQ</b>	3610	40.86	34.04																																			
<b>SciQ</b>	1000	72.50	81.60																																			
<b>SimpleQA</b>	4326	6.80	5.59																																			
<b>Alcuna</b>	2001	-	-																																			
<b>FalseQA</b>	1374	-	-																																			
<b>Self-Aware</b>	1032	-	-																																			
	Table 6: Dataset details. Llama (Llama3-8B-Instruct) and Qwen (Qwen3-8B) are the initial models used in experiments. Acc. (Accuracy) is for reference.																																					
	<b>B Prompts</b>																																					
	During training, we use the following instruction:																																					
	<i>You are a helpful and truthful AI assistant. You should answer the question as briefly as possible, if you don't know, please just say 'I don't know'.</i>																																					
	We use a 6-shot prompt for evaluation across all methods. For the ICL baseline, the prompt consists of six examples of direct answering. In contrast, prompts for abstention-aware methods include a balanced mix of three answering examples and three abstention examples.																																					

### ICL Prompt

Answer the following questions as briefly as possible.

Question: {demo question 1}  
Answer: {demo answer 1}

Question: {demo question 2}  
Answer: {demo answer 2}

...

Question: {input question}  
Answer:

### Abstention-aware Prompt

Answer the following questions as briefly as possible. If you don't know the answer, please simply say "I don't know."

Question: {demo question 1}  
Answer: {demo answer 1}

Question: {demo question 2}  
Answer: I don't know.

...

Question: {input question}  
Answer:

### Uncertainty Prompt

You should answer the question as briefly as possible, then present your confidence. If you are sure about your answer, please say "I am sure" after your answer; otherwise, say "I am unsure".

Question: {demo question 1}  
Answer: {demo answer 1} I am sure.

Question: {demo question 2}  
Answer: {demo answer 2} I am unsure.

...

Question: {input\_question}  
Answer:

### LLM judge Prompt

We are assessing the quality of answers to the following question: *input question*

The following are expected answers to this question: *input ground truth*

The proposed answer is *proposed answer*

Within the context of the question, does the proposed answer mean the same as the expected answer?

Respond only with yes or no.

Here are some examples:

Question: *demo question 1*

Expected answer: *demo ground truth 1*

Proposed answer: *demo correct answer 1*

Response: *yes*

Question: *demo question 2*

Expected answer: *demo ground truth 2*

Proposed answer: *demo wrong answer 2*

Response: *no*

...

Now evaluate the following:

Question: *input question*

Expected answer: *input ground truth*

Proposed answer: *input proposed answer*

Response:

839

840

841

842