

# Learning Auxiliary Tasks Improves Reference-Free Hallucination Detection in Open-Domain Long-Form Generation

Anonymous ACL submission

## Abstract

Hallucination, the generation of factually incorrect information, remains a significant challenge for large language models (LLMs), especially in open-domain long-form generation. Existing approaches for detecting hallucination in long-form tasks either focus on limited domains or rely heavily on external fact-checking tools, which may not always be available.

In this work, we systematically investigate reference-free hallucination detection in open-domain long-form responses. Our findings reveal that internal states (*e.g.*, model’s output probability and entropy) alone are insufficient for reliably (*i.e.*, better than random guessing) distinguishing between factual and hallucinated content. To enhance detection, we explore various existing approaches, including prompting-based methods, probing, and fine-tuning, with fine-tuning proving the most effective. To further improve the accuracy, we introduce a new paradigm, named RATE-FT, that augments fine-tuning with an auxiliary task for the model to jointly learn with the main task of hallucination detection. With extensive experiments and analysis using a variety of model families & datasets, we demonstrate the effectiveness and generalizability of our method, *e.g.*, +3% over general fine-tuning methods on LongFact.

## 1 Introduction

With the recent advancements in model scale and pretraining data, large language models (LLMs) have demonstrated remarkable capabilities in various natural language processing (NLP) tasks (Brown et al., 2020). Despite these successes, hallucination, where models tend to produce content that conflicts with real-world facts, remains a significant challenge (Zhang et al., 2023). Most existing research on hallucination detection has focused on short-form tasks, where the output consists of one or a few tokens. While these methods are effective for short-form content (Manakul

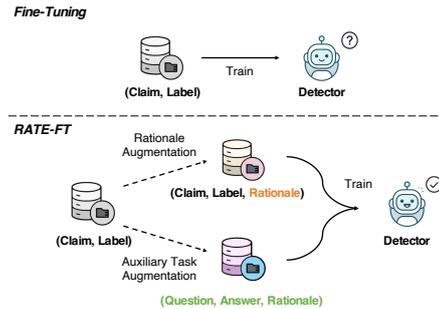


Figure 1: Comparison between Fine-Tuning and RATE-FT for hallucination detection. RATE-FT improves Fine-Tuning by incorporating rationales and an auxiliary task (question answering) into the training process.

et al., 2023; Mahaut et al., 2024; Yehuda et al., 2024; Zhang et al., 2024a), extending them to open-domain long-form generation presents additional complexities and new challenges. Unlike short-form tasks, long-form responses can span hundreds or even thousands of tokens, requiring models to generate detailed and nuanced answers to broad fact-seeking prompts (Wei et al., 2024). This necessitates that LLMs synthesize information across multiple knowledge domains, increasing the risk of generating content that sounds plausible yet is factually incorrect. For example, when answering ‘What is the significance of Amber Room?’, LLMs may generate responses that mix accurate historical information with fabricated details, complicating the task of distinguishing fact from hallucination.

Recent efforts have sought to address hallucination detection in long-form tasks. However, they either focus on limited domains, *e.g.*, biography generation (Min et al., 2023; Fadeeva et al., 2024) or rely heavily on external fact-checking tools or knowledge bases, *e.g.*, Google Search (Wei et al., 2024). While these tools offer valuable support, they are not always available or scalable. This raises an important question: *can we develop hallucination detectors that rely solely on the model itself, without the need for external fact-checking*

resources? So far, little attention has been given to systematically exploring how the model’s own mechanisms can be used for detecting hallucinations in open-domain long-form generation.

To address this gap, we start by investigating hallucination detection in open-domain long-form responses using the model’s internal states, *e.g.*, output probability and entropy. Specifically, we decompose long-form responses into atomized claims using the model and verify each claim’s correctness using Google Search to construct benchmark data following Wei et al. (2024). Our analysis reveals that these internal states alone are insufficient for reliably (*i.e.*, better than random guessing) distinguishing between correct and incorrect claims, indicating that the mechanisms for detecting hallucinations in long-form outputs differ significantly from those in short-form tasks. To enhance detection, we explore several existing methods, including prompting, probing, and fine-tuning LLMs. Our experimental results show that fine-tuning LLMs is the most effective method to detect hallucinations.

Building on this, we introduce a novel method Rationale and Auxiliary Task Enhanced Fine-Tuning (RATE-FT) (Figure 1). Specifically, we convert the original claims into auxiliary question answering (QA) examples for augmentation, providing a complementary learning perspective for the model, which enables better generalization. Additionally, we incorporate collected rationales into the training process for better reasoning. Extensive experiments and analysis using different models demonstrate the effectiveness and generalizability of our approach. Furthermore, we investigate the integration of model uncertainty into hallucination detection in Appendix A.11.

## 2 Are LLMs’ Internal States Sufficient for Open-Domain Long-Form Generation?

The internal states of LLMs, such as output probability and entropy, have been shown to be effective in detecting hallucinations in short-form tasks, where outputs are typically limited to only a few tokens. By analyzing these signals, models can often differentiate between factual and hallucinated information. However, their applicability in open-domain long-form generation remains underexplored. A key question is whether LLMs can depend solely on their internal states to identify hallucinations in long-form generation, without using external fact-checking tools. To answer it,

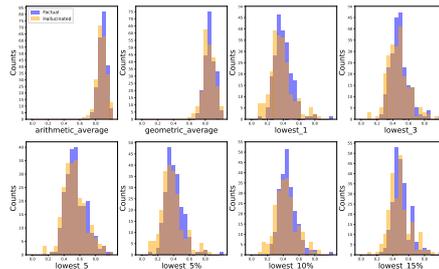


Figure 2: Detection results based on token probability.

we conduct some pilot experiments on LongFact (Wei et al., 2024), a long-form generation dataset spanning 38 different domains. Specifically, for each prompt in the sampled subset (200 prompts), we obtain a long-form response from Llama-3-8B-Instruct with greedy decoding. Following Wei et al. (2024), we employ the model to decompose long-form responses into atomized claims and label them as ‘factual’ or ‘hallucinated’ together with the reasons (see Appendix A.3 for construction details).

For each claim, we mainly focus on two types of internal states to estimate factual confidence following SelfCheckGPT (Manakul et al., 2023): the probability or the entropy (uncertainty) of output tokens. Specifically, we examine the arithmetic and geometric <sup>1</sup> averages of all tokens, the average of tokens with the top- $K$  lowest probability or highest entropy ( $K = 1, 3, 5$ ), and the average of tokens with the top- $P\%$  lowest probability or highest entropy ( $P = 5, 10, 15$ ). The results in Figure 2 and Appendix A.4 suggest that neither internal state reliably, *i.e.*, better than random guessing, predicts the correctness of a given claim, which may be due to the presence of numerous insignificant tokens within the claim, such as stop words. To address this, we consider variants that focus only on output tokens related to entities. The results, shown in Appendix A.4, reveal similar patterns (see Appendix A.5 for a detailed comparison with the findings in Manakul et al. (2023)). We analyze the underlying reasons as follows. In open-domain long-form generation, claims are not limited to a few tokens, which introduces multiple sources of uncertainty. Specifically, *the probability or entropy reflects the model’s confidence in how a claim is expressed, i.e., its confidence in the claim as a sequence of output tokens, rather than in the correctness of the claim.* Different surface forms of the claim yield different confidence levels, leading to unreliable estimates.

Considering the unreliability of LLMs’ internal states in hallucination detection, there are sev-

<sup>1</sup>Commonly known as perplexity

eral promising alternative approaches, including prompting, probing and fine-tuning LLMs, which we explore in the next section.

### 3 Prompting, Probing and Fine-Tuning

Based on a review of the research area, we identify three groups of existing hallucination detection methods, which we discuss below.

**Prompting** Prompting-based approaches involve directly prompting LLMs to assess the correctness of a given claim without additional training. We investigate the following three different methods: (i) Prompting the model to output ‘True’ or ‘False’ for a given claim, referred to as  $\text{Prompt}_{\text{TF}}$ . The probability assigned to the token ‘True’ represents  $P_{\text{factual}}$ , while the probability assigned to ‘False’ represents  $P_{\text{hallucinated}}$ . (ii) Prompting the model to output the *probability* that it considers the given claim to be correct, referred to as  $\text{Prompt}_{\text{Prob}}$ . This number directly represents  $P_{\text{factual}}$ . (iii) SelfCheckGPT, which detects hallucinations by sampling additional responses from the model and assessing inconsistencies between each response and the target claim. The *proportion* of responses that support the claim is taken as  $P_{\text{factual}}$ . Following Manakul et al. (2023), we sample 20 responses for detection.

**Probing** Following Su et al. (2024a), we train a multilayer perceptron (MLP) on the contextualized embeddings of LLMs to perform binary classification for hallucination detection, while keeping the base LLM frozen. The trained MLP outputs  $P_{\text{factual}}$  as an indicator for classification.

**Fine-Tuning** We fine-tune the base LLM with LoRA to enhance its ability to output ‘True’ or ‘False’ for a given claim (Kapoor et al., 2024). Similar to  $\text{Prompt}_{\text{TF}}$ , the probabilities assigned to the tokens ‘True’ and ‘False’ correspond to  $P_{\text{factual}}$  and  $P_{\text{hallucinated}}$ , respectively. Note that LoRA fine-tuning allows us to easily use the original model for general tasks while applying the trained LoRA specifically for hallucination detection.

Following the data construction process outlined in Appendix A.3, we conduct experiments on the full set of LongFact using Llama-3-8B-Instruct. This process yields 2,711 factual and hallucinated claims, which are subsequently split into training (70%), validation (20%), and test (10%) sets. For all three types of methods, we use  $P_{\text{factual}}$  as the classification indicator. Specifically, a claim is classified as ‘factual’ if  $P_{\text{factual}}$  exceeds a predefined threshold; otherwise, it is classified as ‘hal-

Dataset	Method				
	Prompt <sub>TF</sub>	Prompt <sub>Prob</sub>	SelfCheckGPT	Probing	Fine-Tuning
LongFact	69.9	53.4	69.1	74.4	<b>76.1</b>
Biography	72.3	56.3	71.9	77.0	<b>78.2</b>

Table 1: BAcc (%) of existing hallucination detection methods on LongFact and biography generation.

lucinated’. The optimal threshold is determined through a search on the validation set. Consistent with Tang et al. (2024); Chen et al. (2024b), we employ balanced accuracy (BAcc) as the evaluation metric:  $\text{BAcc} = \frac{1}{2}(\frac{\text{TP}}{\text{TP}+\text{FN}} + \frac{\text{TN}}{\text{TN}+\text{FP}})$ , where TP, TN, FP, and FN stand for true/false positives/negatives.

The results of different methods on the test set, as shown in Table 1, indicate that fine-tuning LLMs is the most effective among all existing methods (see Appendix A.6 for an analysis of fine-tuning effectiveness in Out-of-Distribution (OOD) scenarios). While both  $\text{Prompt}_{\text{TF}}$  and SelfCheckGPT achieve decent performance, Probing yields notable improvements by incorporating additional training with labels obtained from external search. Fine-Tuning further enhances performance by updating the internal features of LLMs, enabling more effective learning. In contrast,  $\text{Prompt}_{\text{Prob}}$  performs significantly worse, likely due to LLMs’ tendency to output high probabilities for hallucinated claims, leading to overconfidence. Additionally, we extend the experiments to biography generation (Min et al., 2023). The results presented in Table 1 demonstrate that the observations and conclusions can be generalized to different datasets.

Building on these findings, a natural question arises: can Fine-Tuning be further improved to develop more effective hallucination detectors? We answer this question by *incorporating rationales and an auxiliary task into the training process*.

### 4 Rationale and Auxiliary Task Enhanced Fine-Tuning (RATE-FT)

While hallucination detection is not regarded as a reasoning task in the conventional sense, incorporating Chain-of-Thought (CoT) (Wei et al., 2022) explaining the judgment can still be beneficial for distinguishing factual content from hallucinated information as it enables LLMs to better evaluate the correctness of claims by systematically analyzing underlying components. To examine the impact of rationales, we prompt the model to generate a reasoning path before making a judgment (*i.e.*, ‘True’ or ‘False’), referred to as  $\text{Prompt}_{\text{CoT-TF}}$ . This approach improves performance from 69.9 (using

Dataset	Method				
	Prompt <sub>TF</sub>	Prompt <sub>CoT-TF</sub>	Probing	Fine-Tuning	RATE-FT
LongFact	69.9	74.9	74.4	76.1	<b>79.6</b>
Biography	72.3	74.8	77.0	78.2	<b>80.9</b>

Table 2: BAcc (%) of RATE-FT and baseline methods.

Prompt<sub>TF</sub>) to 74.9, highlighting the effectiveness of incorporating CoT reasoning.

**Augmenting Fine-Tuning with Rationales** Building on the above observation, we augment the fine-tuning dataset with rationales generated by the model during data construction, explaining whether the search results support the claims. Notably, we adopt the ‘label-rationale’ format to maintain the same inference cost as the baseline Fine-Tuning. This allows us to directly derive  $P_{\text{factual}}$  from the first output token without requiring the generation of the complete reasoning path.

Consolidating knowledge through repetition in diverse contexts is a fundamental principle of effective human learning (Ausubel, 2012). For example, medical students deepen their understanding of anatomy by studying diagrams, practicing in simulations, and engaging in hands-on dissections, each offering a unique perspective on the same foundational knowledge. Drawing inspiration from this paradigm, we introduce an auxiliary question answering (QA) task into the fine-tuning process to further strengthen the model’s understanding and enhance its generalization capabilities. This auxiliary QA task serves as a complementary component to the primary hallucination detection task, offering the model an alternative but closely related perspective on the problem (see Appendix A.7 for more analysis on the auxiliary task).

**Augmenting Fine-Tuning with QA Task** Specifically, for each claim, we first prompt the model to generate a question about the key information within it. If the claim is factual, we ask the model to extract the correct answer directly from the claim and provide an explanation, forming a QA example. For hallucinated claims, we leverage the augmented rationale to guide the model in generating an appropriate correct answer along with an explanation. After constructing these QA examples, they are combined with the original data for fine-tuning.

By integrating these two strategies, we propose **Rationale and Auxiliary Task Enhanced Fine-Tuning (RATE-FT)** (Figure 1). RATE-FT requires the model to systematically analyze and explain its judgments and allows the model to benefit from complementary learning perspectives, reinforcing

Model	Method				
	Prompt <sub>TF</sub>	Prompt <sub>CoT-TF</sub>	Probing	Fine-Tuning	RATE-FT
Llama-3.1-70B-Instruct	73.2	76.8	79.4	80.6	<b>83.8</b>
Mistral-7B-Instruct	61.8	64.1	68.4	70.8	<b>73.4</b>
Qwen2.5-7B-Instruct	72.8	75.5	77.0	78.4	<b>81.1</b>

Table 3: Results using different models.

its understanding of claims through diverse yet interconnected tasks. Following the experimental setup described in Section 3, we show the comparison between RATE-FT and baseline approaches in Table 2, which demonstrates the superiority of RATE-FT across different datasets (see Appendix A.8 for an analysis of the effect of additional data augmentation compared to the auxiliary QA task).

#### 4.1 Further Analysis

**Generalization to Different Models** Our experiments and analysis so far use Llama-3-8B-Instruct as the backbone model. To verify whether the performance gain of RATE-FT is consistent across different backbone models, we extend the experiments to Llama-3.1-70B-Instruct (Dubey et al., 2024), Mistral-7B-Instruct (Jiang et al., 2023), and Qwen2.5-7B-Instruct (Yang et al., 2024) on LongFact (see Appendix A.9 for details on data collection). From the results shown in Table 3, we can observe that RATE-FT consistently outperforms baseline approaches across all models, demonstrating its robustness and generalizability to diverse model architectures and scales.

In addition, we provide a summary of the main contributions, related work, ablation studies, results of incorporating uncertainty for hallucination detection, all prompts used in our experiments, and implementation details in Appendix A.1 ~ A.2, A.10 ~ A.15, respectively.

## 5 Conclusion

In this work, we systematically investigate reference-free hallucination detection in open-domain long-form generation. Our study begins with an analysis of the model’s internal states, demonstrating that these states alone cannot reliably detect hallucinations. We then evaluate several existing approaches, including prompting, probing, and fine-tuning, with fine-tuning emerging as the most effective method. Building on these findings, we introduce Rationale and Auxiliary Task Enhanced Fine-Tuning (RATE-FT), a novel approach that leverages rationales and an auxiliary task to achieve significant improvements in detection performance across two datasets and various LLMs.

344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
  
355  
356  
357  
358  
  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
  
374  
375  
376  
377  
378  
379  
380  
381  
  
382  
383  
384  
385  
386  
387  
  
388  
389  
390  
391  
392  
  
393  
394  
395  
396  
397

## Limitations

One limitation of our work is its focus solely on improving the performance of the hallucination detector. A potential improvement could be to explore leveraging the detector’s feedback as a reward signal to guide LLMs to generate more factual responses. Additionally, developing a more comprehensive benchmark for hallucination detection in open-domain long-form generation that covers a broader range of domains would further enhance its applicability.

## References

David Paul Ausubel. 2012. *The acquisition and retention of knowledge: A cognitive view*. Springer Science & Business Media.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Jifan Chen, Grace Kim, Aniruddh Sriram, Greg Durrett, and Eunsol Choi. 2024a. [Complex claim verification with evidence retrieved in the wild](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3569–3587, Mexico City, Mexico. Association for Computational Linguistics.

Lida Chen, Zujie Liang, Xintao Wang, Jiaqing Liang, Yanghua Xiao, Feng Wei, Jinglei Chen, Zhenghong Hao, Bing Han, and Wei Wang. 2024b. [Teaching large language models to express knowledge boundary from their own signals](#). *arXiv preprint arXiv:2406.10881*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.

Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. 2024. [Fact-checking the output](#)

[of large language models via token-level uncertainty quantification](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9367–9385, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023. [RARR: Researching and revising what language models say, using language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508, Toronto, Canada. Association for Computational Linguistics.

Minda Hu, Bowei He, Yufei Wang, Liangyou Li, Chen Ma, and Irwin King. 2024. [Mitigating large language model hallucination with faithful finetuning](#). *arXiv preprint arXiv:2406.11267*.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ACM Transactions on Information Systems*.

Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023. [Towards mitigating LLM hallucination via self reflection](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1827–1843, Singapore. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. [Mistral 7b](#). *arXiv preprint arXiv:2310.06825*.

Cheongwoong Kang and Jaesik Choi. 2023. [Impact of co-occurrence on factual knowledge of large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7721–7735, Singapore. Association for Computational Linguistics.

Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. 2024. [Unfamiliar finetuning examples control how language models hallucinate](#). *arXiv preprint arXiv:2403.05612*.

Sanyam Kapoor, Nate Gruver, Manley Roberts, Katherine Collins, Arka Pal, Umang Bhatt, Adrian Weller, Samuel Dooley, Micah Goldblum, and Andrew Gordon Wilson. 2024. [Large language models must be taught to know what they don’t know](#). *arXiv preprint arXiv:2406.08391*.

Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Wen-tau Yih, and Xilun Chen. 2024. [Flame: Factuality-aware alignment for large language models](#). *arXiv preprint arXiv:2405.01525*.



568	<b>A Appendix</b>		616
569	<b>A.1 Main Contributions</b>		617
570	The main contributions of our work are twofold.		618
571	• We are the first to systematically investigate		619
572	reference-free hallucination detection in open-		620
573	domain long-form generation by analyzing a rep-		621
574	resentative set of existing methods.		
575	• We introduce a novel approach that incorporates		
576	rationales and an auxiliary question answering		
577	task into fine-tuning, achieving significant perfor-		
578	mance improvements.		
579	<b>A.2 Related Work</b>		
580	Large Language Models (LLMs) often generate		
581	content that appears plausible but is factually un-		
582	supported, a phenomenon commonly known as hal-		
583	lucination (Zhang et al., 2023). Based on whether		
584	the hallucinated content contradicts read-world		
585	facts or the input context, hallucination can be cat-		
586	egorized into two main groups: factuality halluci-		
587	nation and faithfulness hallucination (Huang et al.,		
588	2023). Extensive research has been conducted on		
589	exploring the causes (Onoe et al., 2022; Kang and		
590	Choi, 2023; Wei et al., 2023; Liu et al., 2024), de-		
591	tection (Min et al., 2023; Zhao et al., 2023; Chen		
592	et al., 2024a; Fadeeva et al., 2024; Wei et al., 2024),		
593	and mitigation (Gao et al., 2023; Ji et al., 2023; Tian		
594	et al., 2024; Zhang et al., 2024b; Kang et al., 2024;		
595	Lin et al., 2024) of hallucination in LLMs. How-		
596	ever, most existing hallucination detection methods		
597	have primarily focused on short-form tasks, where		
598	the output consists of one or a few tokens. In this		
599	work, we shift the focus to the more challenging		
600	problem of reference-free hallucination detection		
601	in open-domain long-form generation, where out-		
602	puts are substantially longer and require a more		
603	nuanced evaluation of actuality.		
604	<b>A.3 Benchmark Construction Details</b>		
605	For each prompt in the sampled subset (200		
606	prompts), we obtain a long-form response from		
607	Llama-3-8B-Instruct with greedy decoding. Fol-		
608	lowing Wei et al. (2024), we employ the model		
609	to decompose long-form responses into atomized		
610	claims and assess whether each claim is relevant to		
611	answering the corresponding prompt. For each rele-		
612	vant claim, we use the model to generate multi-step		
613	Google Search queries and reason about whether		
614	the search results support the claim. Claims sup-		
615	ported by the search results are labeled as “factual”,		
	while those contradicted by the results are cate-		616
	gorized as “hallucinated”. After construction, we		617
	obtain 2394 factual claims and 223 hallucinated		618
	claims, respectively. We then randomly selected		619
	an equal number (223) of factual and hallucinated		620
	claims for experiments.		621
	<b>A.4 Hallucination Detection Results using</b>		622
	<b>Internal States</b>		623
	We show the hallucination detection results using		624
	different internal states in Figure 3 ~ 5		625
	<b>A.5 Detailed Comparison with Findings in</b>		626
	<b>SelfCheckGPT</b>		627
	(i) While SelfCheckGPT (Manakul et al., 2023)		628
	explores several internal states of LLMs, our work		629
	covers a broader range of variants. As illustrated in		630
	Section 2, we examine the arithmetic and geometric		631
	averages (perplexity) of all tokens, the average of		632
	tokens with the top- $K$ lowest probability or highest		633
	entropy ( $K = 1, 3, 5$ ), and the average of tokens		634
	with the top- $P\%$ lowest probability or highest en-		635
	tropy ( $P = 5, 10, 15$ ). In contrast, SelfCheckGPT		636
	only examines the arithmetic average of all tokens		637
	and the average of tokens with the top-1 lowest		638
	probability or highest entropy.		639
	(ii) Our findings differ significantly from those re-		640
	ported in SelfCheckGPT. While SelfCheckGPT		641
	suggests that LLM probabilities correlate well with		642
	factuality, our experiments demonstrate that nei-		643
	ther internal state reliably, <i>i.e.</i> , better than random		644
	guessing, predicts the correctness of a given claim.		645
	One possible explanation for this is the presence		646
	of many insignificant tokens, such as stop words,		647
	within the claim. To address this, we further in-		648
	vestigate variants that focus only on output tokens		649
	related to entities (Appendix A.4), and the results		650
	exhibit similar patterns. Importantly, our findings		651
	are consistent with those in Kapoor et al. (2024).		652
	<b>A.6 Out-of-Distribution Results</b>		653
	We verify the effectiveness of fine-tuning in Out-		654
	of-Distribution (OOD) scenarios by training the		655
	model on LongFact and evaluating its performance		656
	on Biography. The results reported in Table 4		657
	demonstrate that fine-tuning effectively generalizes		658
	to OOD scenarios.		659
	<b>A.7 More Analysis on Auxiliary Task</b>		660
	<b>Comparison with F2</b> F2 (Hu et al., 2024) also		661
	integrates rationales and auxiliary tasks into the		662

training process. However, its main goal is to enhance the faithfulness of model responses while we focus on improving the accuracy of hallucination detection.

**Further Clarification on Motivation** The underlying motivation for introducing the auxiliary question answering (QA) task into fine-tuning is that hallucination detection and mitigation are complementary and closely related tasks. This auxiliary QA task—where a question about the key information in the claim is posed, and the model is trained to provide the correct answer—helps improve the factuality of the model’s responses through supervised fine-tuning. It acts as a complementary component to the primary hallucination detection task, offering the model an alternative yet closely related perspective, thereby enhancing its generalization capabilities.

### A.8 Additional Data Augmentation versus Auxiliary QA Task

To isolate the effect of additional data augmentation versus the auxiliary QA task, we design two variants: (i) we paraphrase the original claim using GPT-4 for data augmentation and fine-tune the model on the combined data, referred to as Fine-Tuning<sub>para</sub>, which has roughly the same amount of training data as RATE-FT; and (ii) we reduce the training data for RATE-FT by half (approximately the same amount as Fine-Tuning), referred to as RATE-FT<sub>half</sub>. We conduct experiments on LongFact using Llama-3-8B-Instruct and present the results in Table 5 and 6, which demonstrate that the performance improvement primarily comes from our designed auxiliary task, rather than from additional data augmentation.

### A.9 Data Collection Process for Other Models

When conducting experiments using other models, we follow the exact same settings as those used for Llama-3-8B-Instruct. Specifically, for each prompt, we obtain a long-form response from the model under investigation with greedy decoding. Following Wei et al. (2024), we employ the model to decompose long-form responses into atomized claims and assess whether each claim is relevant to answering the corresponding prompt. For each relevant claim, we use the model to generate multi-step Google Search queries and reason about whether the search results support the claim. Claims supported by the search results are labeled as “factual”, while those

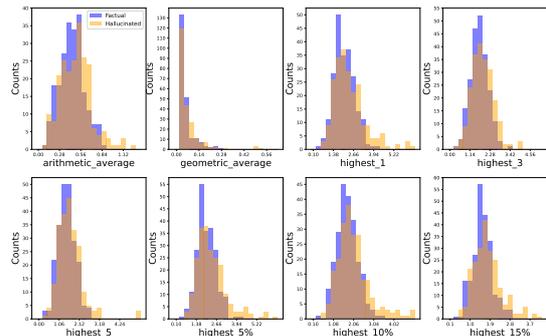


Figure 3: Hallucination detection results based on token entropy (uncertainty).

Prompt <sub>TF</sub>	Prompt <sub>Prob</sub>	SelfCheckGPT	Probing	Fine-Tuning
72.3	56.3	71.9	71.1	<b>74.7</b>

Table 4: Results of different methods in OOD scenarios.

contradicted by the results are categorized as “hallucinated”.

Our constructed benchmarks align well with Su et al. (2024b), as both include responses and internal states from various LLMs. The key difference is that the LLMs we investigate are all modern models (Llama-3-8B-Instruct, Llama-3.1-70B-Instruct, Mistral-7B-Instruct, and Qwen2.5-7B-Instruct), whereas the models used in Su et al. (2024b) are relatively outdated (such as LLaMA-2 and GPT-J).

### A.10 Ablation Study

We analyze the contribution of different components of RATE-FT by investigating the variant of RATE-FT without the auxiliary task (*w.o.* aux). Table 7 presents the performance of different methods, highlighting that each component plays an important role in achieving the overall performance.

### A.11 Incorporating Uncertainty for Hallucination Detection

To enhance hallucination detection, we propose incorporating model uncertainty into the detection process, enabling a hybrid pipeline that combines the strengths of the model and external tools. Specifically, when the model is uncertain about whether a claim is factual or hallucinated, we leverage external tools to handle ambiguous cases, improving overall performance. The process involves setting two thresholds,  $\alpha_{low}$  and  $\alpha_{high}$ , for classification. A claim is classified as ‘factual’ if  $P_{factual} > \alpha_{high}$  and ‘hallucinated’ if  $P_{factual} < \alpha_{low}$ . Claims falling between these thresholds are classified as

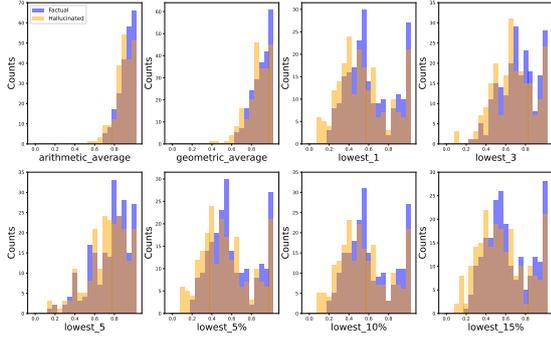


Figure 4: Hallucination detection results based on the probability of entity-related tokens.

Fine-Tuning <sub>para</sub>	RATE-FT
76.8	<b>79.6</b>

Table 5: Comparison between Fine-Tuning<sub>para</sub> and RATE-FT.

Fine-Tuning	RATE-FT <sub>half</sub>
76.1	<b>78.5</b>

Table 6: Comparison between Fine-Tuning and RATE-FT<sub>half</sub>.

‘unknown’ and delegated to external tools for further evaluation. Assuming the external tools’ output is the ground truth, predictions classified as ‘unknown’ are treated as correct. To evaluate the hybrid pipeline, we define the BAcc-unknown metric as follows:

$$\text{BAcc-unknown} = \frac{1}{2} \left( \frac{\# \text{ Correct Factual Predictions}}{\# \text{ Total Factual Claims}} + \frac{\# \text{ Correct Hallucinated Predictions}}{\# \text{ Total Hallucinated Claims}} \right) \quad (1)$$

The optimal thresholds,  $\alpha_{low}$  and  $\alpha_{high}$ , are determined through a search on the validation set. This process ensures that BAcc on the validation set exceeds 70%, while also maximizing BAcc-unknown. The goal is to strike a balance between performance and efficiency by achieving high BAcc-unknown without generating an excessive number of ‘unknown’ predictions, which could substantially increase detection costs. We conduct experiments on LongFact using Llama-3-8B-Instruct and report the results in Table 8, which demonstrate that incorporating model uncertainty greatly enhances hallucination detection, as evidenced by the BAcc-unknown metric’s superior performance compared to standard BAcc in resolving ambiguous cases. Moreover, RATE-FT continues to outperform all

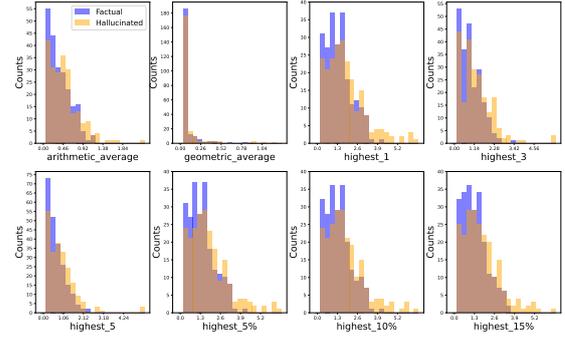


Figure 5: Hallucination detection results based on the entropy of entity-related tokens.

Dataset	Method		
	Fine-Tuning	w.o. aux	RATE-FT
LongFact	76.1	77.5	<b>79.6</b>
Biography	78.2	79.4	<b>80.9</b>

Table 7: Results of different ablations.

Prompt <sub>CoT-TF</sub>	Probing	Fine-Tuning	RATE-FT
80.4	81.1	82.4	<b>85.0</b>

Table 8: BAcc-unknown (%) of different methods on Longfact with Llama-3-8B-Instruct.

other methods with respect to the BAcc-unknown metric, highlighting its robustness and effectiveness.

## A.12 Prompt for Output Extraction

After decomposition, the atomized claims may differ from the original expression in the response. To address this, we use the prompt shown in Figure 6 to retrieve the original output corresponding to a given atomized claim.

## A.13 Prompts for Baseline Approaches

Figure 7 illustrates the prompts used for different prompting methods. The prompt used for constructing training data in Probing and Fine-Tuning is the same as the prompt employed by the Prompt<sub>TF</sub> method.

## A.14 Prompts Used in RATE-FT

Figure 8 presents all the prompts used in RATE-FT.

## A.15 Implementation Details

For Prompt<sub>TF</sub> and Prompt<sub>Prob</sub>, we obtain the response from the model with greedy decoding. Following Manakul et al. (2023), we set the temperature to 1.0 and generate 20 additional responses for

**Prompt**

Your task is to extract the original text corresponding to the given claim from the original response. When presented with an original response and a claim, reply with the original text. Make sure that your response is exactly the same as the original text and enclosed in 'boxed'.

Original response: {response}  
 Claim: {claim}

Figure 6: Prompt for extracting the original output given an atomized claim.

**Prompt<sub>tr</sub>**

Your task is to determine the correctness of the given claim. When presented with a claim, reply with 'True' or 'False'. Make sure that your response is exactly 'True' or 'False' without any extra commentary whatsoever.

Claim: {claim}  
 Response:

**Prompt<sub>tr+ex</sub>**

Your task is to determine the correctness of the given claim. When presented with a claim, first explain the solution and then enclose the ultimate answer ('True' or 'False') in 'boxed'.

Claim: {claim}  
 Response:

**Prompt<sub>prob</sub>**

Your task is to provide the probability that the given claim is correct. When presented with a claim, reply with a number between 0.0 and 1.0. Make sure that your response is exactly a number between 0.0 and 1.0 without any extra commentary whatsoever.

Claim: The sun rises in the east and sets in the west.  
 Response: 1.0

Claim: Humans have four arms and three heads.  
 Response: 0.0

Claim: The human nose can detect over 1 trillion different scents.  
 Response: 0.82

Claim: The next president of South Korea will be a woman.  
 Response: 0.29

Claim: {claim}  
 Response:

**SelfCheckGPT**

Context: {context}

Sentence: {sentence}

Is the sentence supported by the context above? Answer: Yes or No.

Answer:

Figure 7: Prompts for different prompting methods.

### SelfCheckGPT.

We evaluate 4 different types of contextualized embeddings for Probing: (1) the final token from the last layer ( $type_1$ ), (2) the average of all tokens in the last layer ( $type_2$ ), (3) the average of the final token across all layers ( $type_3$ ), and (4) the average of  $type_1$  and  $type_2$  ( $type_4$ ). The optimal embedding type, along with other hyperparameters, *e.g.*, learning rate, is selected through a search on the validation set. For Fine-Tuning and RATE-FT, we leverage the LLaMA-Factory library (Zheng et al., 2024) and perform a search on the validation set for important hyperparameters.

**Prompt for 'label-rationale' Format**

Your task is to determine the correctness of the given claim. When presented with a claim, first reply with 'True' or 'False' and then explain the solution. Make sure that your response starts with 'True' or 'False'.

Claim: {claim}  
 Response: {True/False}, {explanation}

**Prompt for Question Answering**

Answer the following question and provide the explanation.

Question: {question}  
 Answer: {answer}  
 Explanation: {explanation}

**Prompt for Question Generation (Correct Claim)**

Given a correct claim and why it is correct, first identify the key information in the claim, then transform it into a question and a correct answer (keep the answer as concise as possible) about the key information, finally give the explanation (keep it different from the given reason). Make sure that your response follows the format 'Question: {question}|Correct answer: {correct answer}|Explanation: {explanation}'.

Correct claim: {correct claim}  
 Reason: {reason}  
 Response:

**Prompt for Question Generation (Wrong Claim)**

Given a wrong claim and why it is wrong, first identify the key information in the claim, then transform it into a question and a correct answer (keep the answer as concise as possible) about the key information, finally give the explanation (keep it different from the given reason). Make sure that your response follows the format 'Question: {question}|Correct answer: {correct answer}|Explanation: {explanation}'.

Wrong claim: {wrong claim}  
 Reason: {reason}  
 Response:

Figure 8: Prompts for different components of RATE-FT.

789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801