



PRISM: Probing Reasoning, Instruction, and Source Memory in LLM Hallucinations

Yuhe Wu¹, Guangyu Wang^{2,3}, Yuran Chen³, Jiatong Zhang³, Yutong Zhang³,
Yujie Chen⁴, Jiaming Shang⁵, Guang Zhang^{1*}, Zhuang Liu^{3*}

¹HKUST(GZ) ²NYUSH ³DUFE
⁴CUHK(SZ) ⁵CUFE

✉: guangzhang@hkust-gz.edu.cn, liuzhuang@dufe.edu.cn 🌐: <https://acl-prism.cc>

Abstract

As large language models (LLMs) evolve from conversational assistants into agents capable of handling complex tasks, they are increasingly deployed in high-risk domains. However, existing benchmarks largely rely on mixed queries and posterior evaluation, output-level scoring, which quantifies hallucination severity but offers limited insight into *where* and *why* hallucinations arise in the generation pipeline. We therefore reformulate hallucination evaluation as a diagnostic problem and propose **PRISM**, a controlled benchmark that disentangles hallucinations into four dimensions: knowledge missing, knowledge errors, reasoning errors, and instruction-following errors, grounded in three stages of generation (memory, instruction, and reasoning). **PRISM** contains 9,448 instances across 65 tasks and supports fine-grained, stage-aware diagnostic evaluation. Evaluating 24 mainstream open-source and proprietary LLMs, we uncover consistent trade-offs across instruction following, memory retrieval, and logical reasoning, showing that mitigation strategies often improve specific dimensions at the expense of others. We hope **PRISM** provides a framework for understanding the specific mechanisms behind LLMs hallucinations, ultimately accelerating the development of trustworthy large language models.

1 Introduction

LLMs have become capable of handling complex tasks (Wang et al., 2024b; Zhang et al., 2025a; Liu et al., 2025; Xi et al., 2025), facilitating their application in high-risk domains such as medical diagnosis (Singhal et al., 2023; Thirunavukkarasu et al., 2023), legal consulting (Guha et al., 2023; Cui et al., 2024), and scientific discovery (Boiko et al., 2023; Bran et al., 2024). While current models perform well on general benchmarks (Hendrycks et al., 2021; Achiam et al., 2024), they frequently

generate factually inconsistent content (Alansari and Luqman, 2026) when encountering outdated concepts (Kandpal et al., 2023; Mallen et al., 2023), dynamic information (Kasai et al., 2023; Vu et al., 2024), or complex reasoning and instruction constraints (Dziri et al., 2023; Lanham et al., 2023; Heyman and Zylberberg, 2025). Such unfaithfulness not only erodes user trust but also constitutes potential safety hazards in critical decision-making scenarios (Thirunavukkarasu et al., 2023; Zhang et al., 2025b). Consequently, the evaluation of hallucinations has emerged as a fundamental challenge that the research community needs to overcome.

Despite the growing interest in quantifying hallucinations (Lin et al., 2022; Li et al., 2023), existing benchmarks have clear limitations in answering the fundamental question of why models fail.

First, current benchmarks often mix different queries, which prevents us from testing skills in isolation. As shown in Figure 1 (Left Top), benchmarks like TruthfulQA (Lin et al., 2022), HaluEval (Li et al., 2023), and FreshQA (Vu et al., 2024) typically use mixed queries. When a model fails on these, the reason is ambiguous: did it fail to retrieve the right data, make a logical error, or simply ignore the instructions? Second, most evaluations focus only on the final output. Even detailed methods like FActScore (Min et al., 2023), and HALoGEN (Ravichander et al., 2025) depend on posterior evaluation after generation. Relying on such outcome assessments introduces unavoidable bias from both human and model evaluators. Furthermore, this approach fails to test specific inputs to isolate the error. Without knowing exactly where the process broke, fixing the model is much harder. An abstract comparison is shown in Table 1.

As illustrated in Figure 1 (Right), improvement strategies for different mechanisms often involve inherent trade-offs: for instance, strong instruction fine-tuning to fix formatting errors can acciden-

*Corresponding authors.

tally hurt rigorous reasoning capabilities (Ouyang et al., 2022; Peng et al., 2023), while indiscriminate knowledge injection may cause catastrophic forgetting (Zhai et al., 2023). Consequently, to achieve reasonable optimization, we must address a fundamental question:

❖ *How can we establish a framework to pinpoint failures in memory, reasoning, and instruction following, thereby guiding optimization for hallucination mitigation?*

To address the aforementioned challenges and establish a trustworthy diagnostic framework, we propose *PRISM*, an evaluation benchmark grounded in the interactive pipeline of LLMs. Based on the three generation stages of instruction following, memory retrieval, and reasoning, we exactly categorize hallucination phenomena into four independent failure dimensions:

- **Knowledge Error (KE):** The model’s parametric knowledge stores incorrect or outdated information.
- **Knowledge Missing (KM):** The model’s parametric knowledge lacks the correct information required to answer the question.
- **Reasoning Error (RE):** The model possesses the necessary facts but fails to combine them through logic or reasoning.
- **Instruction Following Error (IFE):** The model possesses correct knowledge and reasoning capabilities, but its output violates explicit constraints provided by the user.

This design enables us to locate specific shortfalls of the model in the generation stages. Our main contributions are summarized as follows:

- We propose a cognitive pipeline failure framework that defines hallucinations as dimensions in KE, KM, RE, and IFE. We then build *PRISM*, a benchmark of 9448 samples that isolates these factors and pinpoints model weaknesses for reproducible analysis.
- We conducted a comprehensive evaluation of 24 proprietary and open-source LLMs across 4 dimensions and 65 sub-tasks to assess the causes of hallucinations in different model types and encourage the training of hallucination-specific LLMs.

- Building on *PRISM*, we examined the performance trade-offs of common hallucination mitigation strategies. Furthermore, we constructed a toy dataset to support case-based empirical studies, allowing us to reveal the internal mechanisms of LLMs during KE and KM memory-based issues and analyze the relationship between IFE and LLMs efficiency. These findings guide the design of balanced mitigation strategies.

2 Benchmark Construction

To achieve precise attribution of hallucination mechanisms, our data construction follows the principle of orthogonality, meaning that each data subset aims to independently test a single failure mode to the greatest extent possible.

2.1 Data Source

To achieve attribution of hallucination dimensions, *PRISM*’s data construction strictly follows the orthogonality principle: each subset is designed to test a single failure mode. As shown in Table 1, existing benchmarks suffer from limited evaluation scope and a lack of variable control, making metrics incapable of revealing the causes of errors. Consequently, we construct a corpus that strictly partitions data into parametric knowledge-dependent and Reasoning and instruction-dependent categories, enabling the isolated probing of specific hallucination modes. Further detailed source lists are available in Appendix C.

Sources for Parametric Knowledge Tasks. This category aims to define the accuracy and boundaries of the model’s internal memory by comparing it against external objective facts. We collected two types of raw corpus:

- **Factual Data:** To ensure the factual consistency of our evaluation standards, we selected *Wikipedia* and *Baidu Baike* as the primary sources for KE tasks. Compared to unfiltered web texts, these corpus have lower noise. We focused on collecting long-tail and ambiguous entries that are not in the LLM’s source memory base, used to test the model’s memory coverage and the ability to distinguish specific entities.
- **Out-of-Distribution (OOD) Data:** To evaluate the model’s ability to identify unknown information, we first established a temporal news corpus

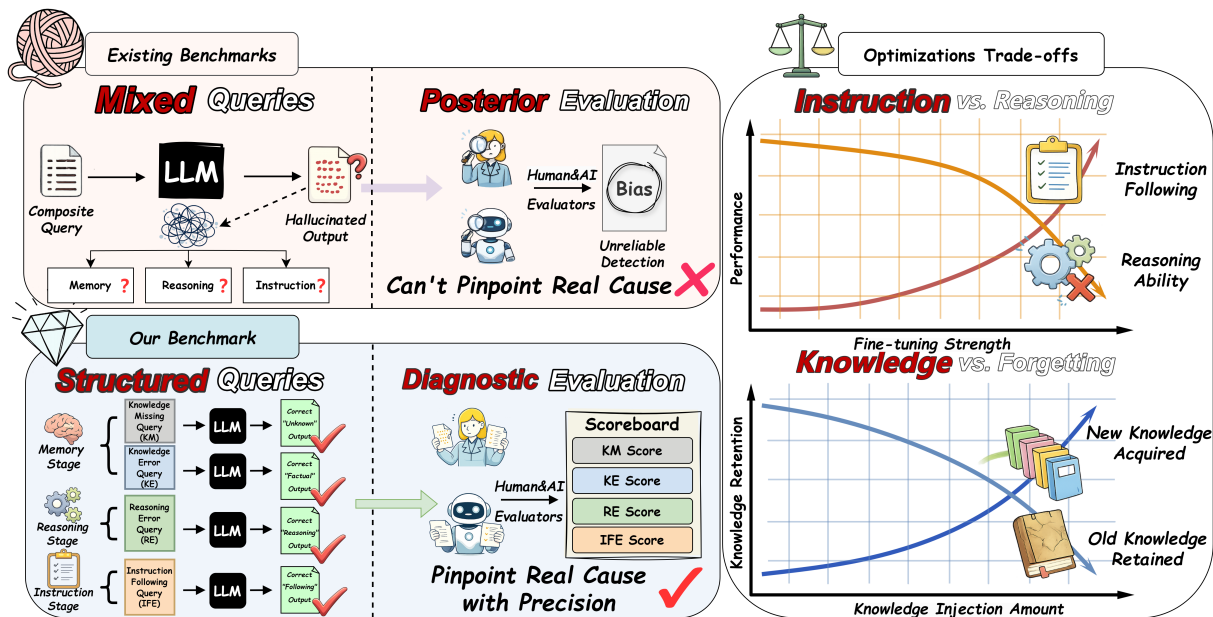


Figure 1: Overview of the **PRISM** framework and optimization trade-offs. The left panel contrasts the mixed query design of existing benchmarks with our structured approach that isolates cognitive stages to pinpoint failure dimensions like KE, KM, RE, and IFE. The right panel illustrates performance trade-offs where enhancing instruction following compromises reasoning ability and knowledge injection leads to the forgetting of retained information.

by collecting news reports and paper abstracts published by CNN, Reuters, and arXiv between March 2025 and November 2025. As this period postdates the training cutoff of most baseline models, these materials constitute a test source for future information. Second, we constructed a fictional entity corpus. Rather than being collected from the real world, this content consists of generated counterfactual descriptions created by setting specific attributes. Finally, to cover information that is real but non-public, we introduced private domain data.

Sources for Reasoning and Instruction Tasks.

This type of data is designed to reduce the model’s reliance on parametric knowledge and to focus on evaluating its ability to perform logical reasoning and execute rules under a given context. We built two corpus to keep the focus on the reasoning process rather than memory retrieval.

- **Self-Contained Reasoning Data:** To ensure the reasoning process is isolated from external knowledge noise, we prioritized task types where the solution premises are strictly embedded within the input context. In addition to introducing competition problems, such as the IMO, to cover formal logic and mathematical proofs, we also incorporated code generation tasks. Given their deterministic execution logic,

they provide the most ideal ambiguity-free reasoning environment for the model.

- **Complex Instruction Data:** Beyond covering everyday basic instructions, we specifically construct a high-constraint corpus, and we use automated templates to generate adversarial synthetic data. This corpus simulates real scenarios where instruction violations happen due to competition for attention resources under multi-dimensional stacked constraints, including negative semantics that forbid specific words, format locking that requires strict JSON output, and various limits on length and language.

2.2 Construction Pipeline

To construct **PRISM**, we design a three-stage pipeline as illustrated in Figure 2.

Data Collection. In this initial stage, a corpus is collected from authoritative sources and cleaned via noise removal.

Multi-agent Data Construction. Next, we adopt a multi-agent framework to construct data: (i) schema normalizer agent; (ii) evidence retriever agent; (iii) type classifier agent; (iv) quality scoring agent, enabling each agent to focus on a specific step for improving the results, thereby enhancing both the efficiency and quality of the data.

Benchmark	Evaluation Scope				Methodological Design	
	KE	KM	RE	IFE	Variable Control	Diag. Mode
TruthfulQA (Lin et al., 2022)	✓	✗	✗	✗	✗	✗
HaluEval (Li et al., 2023)	✓	✓	✓	✗	✗	✗
FActScore (Min et al., 2023)	✓	✓	✗	✗	✗	✗
FELM (Chen et al., 2023)	✓	✗	✓	✗	✗	✗
FreshQA (Vu et al., 2024)	✓	✓	✗	✗	✓	✗
FollowBench (Jiang et al., 2024b)	✗	✗	✓	✓	✗	✗
HALoGEN (Ravichander et al., 2025)	✓	✓	✓	✗	✗	✗
HalluLens (Bang et al., 2025)	✓	✓	✓	✗	✗	✗
PRISM (Ours)	✓	✓	✓	✓	✓	✓

Table 1: Comparison of hallucination evaluation benchmarks. *PRISM* uniquely achieves *comprehensive evaluation scope*, *causative explainability*, and *decoupled probing-based diagnosis*.

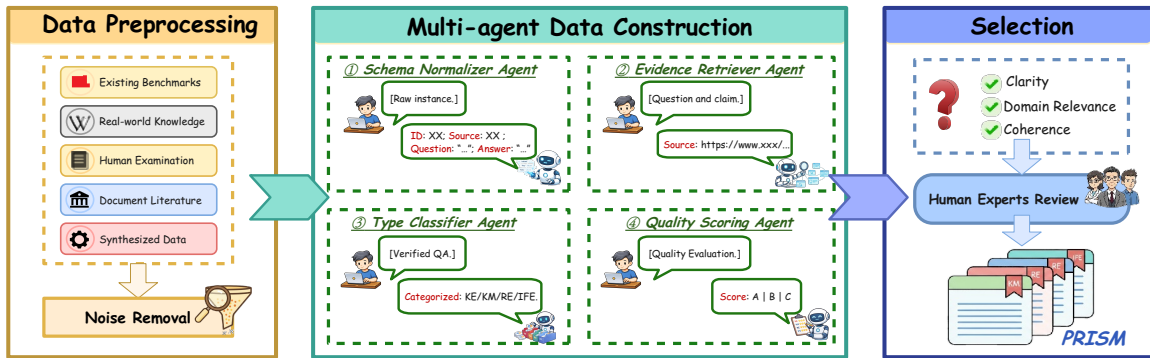


Figure 2: The Three-phase Pipeline of *PRISM* Benchmark Construction

Human Selection. Domain experts select instances for clarity, relevance, and coherence to curate the *PRISM*. Detailed construction procedures are provided in Appendix D.

2.3 Data Statistics

PRISM contains a total of 9,448 evaluation instances, covering 4 failure dimensions and 65 specific sub-tasks. Among them, 2,995 are RE samples (31.7%), 2,442 are IFE samples (25.9%), 2,078 are KM samples (22.0%), and 1,933 are KE samples (20.5%).

The distribution of the data is illustrated in Figure 3. The left side of the figure displays a sunburst chart where the inner circle represents the four primary failure dimensions, and the outer ring corresponds to sub-task indices ranging from 1 to 65. The right side of the figure lists the detailed mapping for each index, specifying the category name and the exact sample count for each sub-task.

3 Experiments

PRISM is designed to evaluate the reliability of LLMs and provide guidance for model optimization. To meet these objectives, we structure our experiments around three research questions that establish performance baselines, explain underlying causes, and identify pathways for improvement:

- **RQ1:** What is the overall performance of LLMs on *PRISM*, and how do different error types vary across model families and scales?
- **RQ2:** How effective are common mitigation strategies, and which methods best address specific types of hallucinations?
- **RQ3:** Why do different types of hallucinations show consistent patterns in the internal representations of LLMs?

To provide a solid foundation for investigating these questions, we begin by outlining the experimental setup, focusing on model selection and evaluation protocols.

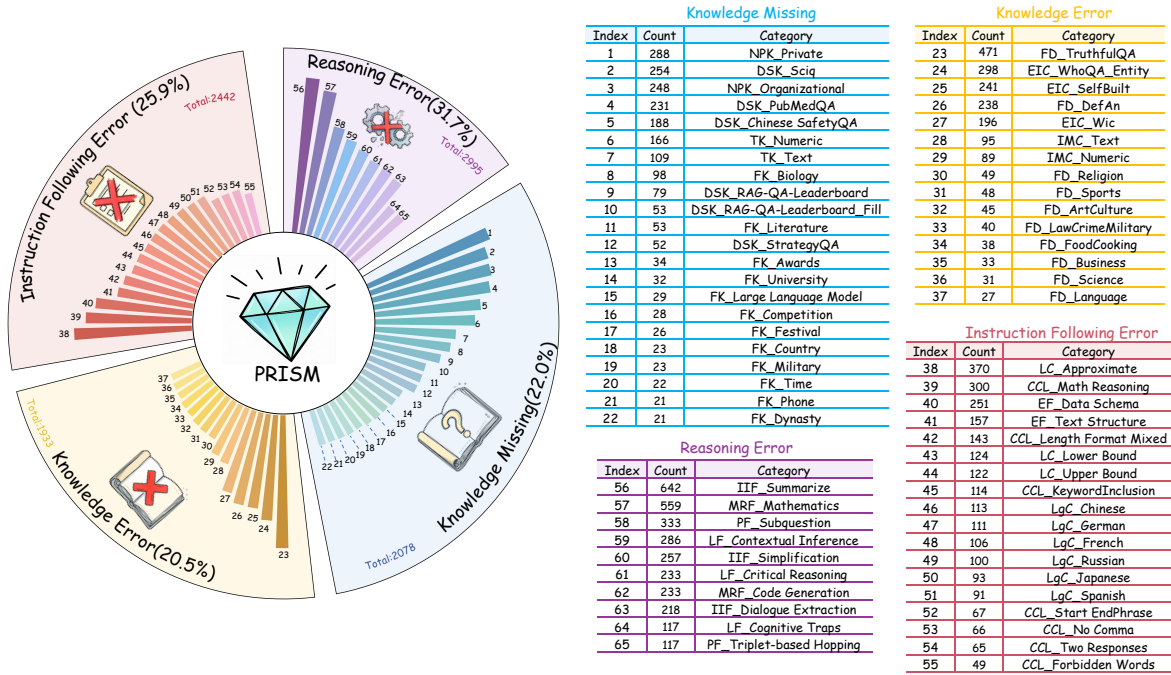


Figure 3: The hierarchical distribution of **PRISM**. The inner circle represents the four primary failure dimensions, while the outer ring details 65 sub-tasks. For consistency, we define the abbreviations as follows: DSK = Domain-Specific Knowledge, FK = Fictional Knowledge, TK = Timely Knowledge, NPK = Non-Public Knowledge, FD = Factual Distortion, IMC = Intra-Memory Conflict, EIC = Entity-Identity Confusion, LF = Logical Fallacy, PF = Procedural Failure, IIF = Information Integration Failure, MRF = Mathematical Reasoning Failure, EF = Explicit Format, LC = Length Constraints, LgC = Language Constraints, and CCL = Complex & Cognitive Load.

3.1 Evaluation Setup

Model Selection. We evaluate 24 representative LLMs under the few-shot setting. To ensure comprehensive coverage, the evaluated LLM series includes both open-source and proprietary models, such as GPT (OpenAI, 2024), Gemini (Comanici et al., 2025), Llama (AIMeta, 2025), Claude (Anthropic, 2025), DeepSeek (Guo et al., 2025), GLM (Team et al., 2024), Qwen (Team, 2024), and Grok (xAI, 2025).

Evaluation Metrics. We employ distinct metrics for each subtask to enable a hallucination comparison.

- **Accuracy:** For closed-ended tasks, we employ standard Accuracy:

$$\text{Acc} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[\hat{y}_i = y_i]$$

- **LLM-Eval:** For open-ended tasks, we adopt a LLM evaluator following LLM-EVAL (Lin and Chen, 2023; Zheng et al., 2023), which produces a scalar score $s \in [0, 5]$. The

prompts and the consistency evaluation between the LLM-as-a-judge and human annotations are provided in Appendix J.

- **Hallucination Rate:** We first map all task metrics to a unified percentage score $S \in [0, 100]$:

$$S = \begin{cases} 100 \cdot \text{Acc}, & \text{for closed-ended tasks,} \\ 100 \cdot \frac{s}{5}, & \text{for open-ended tasks,} \end{cases}$$

and define the hallucination rate as its complement:

$$\mathcal{H} = 100 - S.$$

- **\mathcal{H} -Score:** Let \mathcal{H}_d denote the macro-averaged hallucination rate for each dimension $d \in D = \{\text{KE, KM, RE, IFE}\}$. We define

$$\mathcal{H} - \text{Score} = \frac{1}{4} \sum_{d \in D} \mathcal{H}_d.$$

Sampling Parameters. To ensure fair and comparable evaluations across diverse models, we control generation randomness using temperature and top p sampling (Holtzman et al., 2020; OpenAI,

2025), exploring temperature values in {0, 0.2, 0.4, 0.6, 0.8} and top-p values in {0.6, 0.8, 0.9, 0.95}. For closed-ended tasks, we conducted the grid search on the DeepSeek-R1-Distill-32B, identifying temperature = 0.8 and top-p = 0.8 as optimal for the KM, KE, and IFE subsets, while adopting temperature = 0.4 and top-p = 0.95 for RE to enhance reasoning stability. For open-ended tasks, we performed the grid search on GPT-4o, yielding temperature = 0.8 and top-p = 0.8 for peak average performance. The full grid results are reported in Appendix G.

3.2 Evaluation Results (RQ1)

Table 2 comprehensively reports the hallucination rates across four core evaluation dimensions along with the composite \mathcal{H} -Score metrics. Claude-Opus-4.5, Gemini-3-Pro, and Gemini-3-Flash achieve the top three rankings with the lowest \mathcal{H} -Scores (13.90%, 14.29%, and 15.55%, respectively), indicating more robust and balanced performance across diverse error types. In terms of dimensional distribution, RE and KE exhibit significant performance gaps among models. In contrast, error rates in the KM dimension remain relatively low for most models, rendering this dimension less of a dominant factor in determining the final ranking. The analysis in Figure 4 further reveals only partial consistency across dimensions: KE shows the strongest correlation with RE and the weakest with IFE. More detailed experimental results are provided in Appendix F.

In the comparison of model types, open-source models have gradually approached proprietary models in the KM dimension but still lag significantly in RE and IFE, particularly on complex tasks requiring the integration of multiple clues and consistent reasoning. Furthermore, experiments indicate that adopting step-by-step explanation strategies does not yield consistent performance gains. This is likely because the hallucination phenomenon is complex, encompassing factual errors, inconsistent reasoning logic, misuse of evidence, and deviation from task requirements during complex instruction following (Ji et al., 2023a; Huang et al., 2025). Additionally, the generated step-by-step explanations may appear plausible on the surface but do not faithfully reflect the internal information upon which the model actually relies during generation (Turpin et al., 2023).

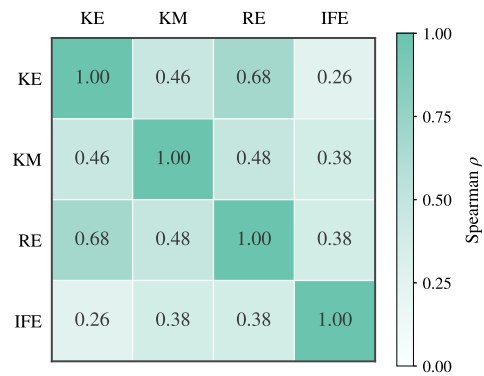


Figure 4: Spearman Correlation of Model rankings cross 4 Dimensions

4 Discussion

Our experimental results reveal a complex relationship between model parameter scale, reasoning capability, and performance. This section first addresses RQ2 by analyzing how different mitigation strategies affect model capabilities. Building on this analysis, we then turn to RQ3 to investigate the potential causes of different hallucination patterns.

4.1 Impact of Mitigation Strategies (RQ2)

To answer RQ2, we examine how hallucination mitigation strategies affect model capability along three dimensions. The first dimension is in-context learning (ICL), where we vary the number of demonstrations to adjust the strength of contextual guidance. The second dimension is instruction tuning, where we introduce Llama-3.1-8B-Instruct to assess behavioral differences induced by general alignment. The third dimension is reasoning tuning, implemented as supervised fine-tuning (SFT) on the reasoning dataset¹, to study how reasoning enhancement influences different error types. To control variables, the latter two settings are compared under the 1-shot configuration. This design aims to reveal whether optimizations targeting instruction following or reasoning trigger cross-dimensional capability tradeoffs.

Table 3 further shows that none of the three strategies delivers stable improvements across dimensions. Instead, the results form a tradeoff structure. Adjusting shots in ICL barely changes the overall conclusion. The 0-shot setting slightly degrades performance, while the 3-shot setting brings

¹Trained on the gsm8k-reasoning dataset. Available at: <https://huggingface.co/datasets/thesven/gsm8k-reasoning>

Type	Model	Size	Think	KE	KM	RE	IFE	\mathcal{H} -Score
				(Lower the Better)				
Proprietary LLMs	GPT-5.2	-	✓	17.83%	5.57%	24.67%	15.35%	15.85%
	GPT-5.1	-	✓	25.83%	11.46%	26.01%	25.62%	22.23%
	GPT-4o-20241120	-	✗	19.36%	8.90%	29.72%	16.78%	18.69%
	Gemini-3-Pro	-	✓	16.31%	12.80%	17.19%	10.85%	14.29%
	Gemini-3-Flash	-	✓	16.48%	8.12%	21.59%	16.00%	15.55%
	Gemini-2.5-Pro	-	✓	25.33%	8.99%	21.63%	13.84%	17.45%
	Gemini-2.5-Flash	-	✓	19.29%	8.74%	31.58%	15.51%	18.78%
	Claude-Opus-4.5	-	✓	13.80%	6.35%	19.68%	15.77%	13.90%
	Claude-Sonnet-4.5	-	✓	16.11%	7.03%	23.53%	16.87%	15.89%
	Claude-Haiku-4.5	-	✓	18.13%	7.24%	25.59%	18.15%	17.28%
	Grok-4.1	-	✗	19.20%	17.03%	34.94%	20.04%	22.80%
Grok-4-0709	-	✗	18.35%	14.57%	30.97%	15.74%	19.91%	
Open-source LLMs	DeepSeek-V3.2	685B	✗	17.99%	7.94%	28.31%	14.31%	17.14%
	DeepSeek-R1	671B	✓	18.87%	10.13%	28.04%	11.97%	17.25%
	DeepSeek-R1-Distill-32B	32B	✓	18.14%	11.27%	30.40%	18.71%	19.63%
	Qwen3-235B-Instruct	235B	✓	17.37%	7.00%	25.89%	17.87%	17.03%
	Qwen2.5-72B-Instruct	72B	✗	18.56%	8.61%	27.68%	18.43%	18.32%
	GLM-4.5	355B	✓	18.85%	10.70%	29.34%	15.06%	18.49%
	GLM-4	32B	✗	26.20%	16.17%	35.57%	25.34%	25.82%
	Llama-4-Scout	17B	✗	23.67%	7.58%	32.19%	14.17%	19.40%
	Llama-3.3-70B-Instruct	70B	✗	19.75%	6.24%	29.47%	13.13%	17.15%
	Llama-3.1-8B-Instruct	8B	✗	25.53%	14.04%	54.50%	23.36%	29.36%
	Llama-3-70B-8192	70B	✗	20.70%	7.43%	37.50%	15.83%	20.37%
Llama-3-8B-Instruct	8B	✗	23.15%	18.42%	54.02%	32.53%	32.03%	

Table 2: **Main Results on Hallucination Rates.** All values are hallucination rates \mathcal{H} (%) aggregated over 4 dimensions, together with \mathcal{H} -Score. Models are highlighted with **Best** and **Second Best** within each group.

marginal positive changes, suggesting that additional demonstrations mainly aid format rather than addressing root causes such as knowledge gaps. The Instruct model reduces hallucinations more consistently in knowledge conflict and instruction following dimensions, but it introduces a certain level of loss in the reasoning dimension. By contrast, the Reasoning model achieves the strongest improvements in the reasoning dimension, while showing clear capability degradation in knowledge-related and instruction-related dimensions.

Model	Dataset	Base	ICL		Instruct. Reason.	
		1-shot	0-shot	3-shot	1-shot	1-shot
Llama3.1-8B	KE	43.15	43.88 (+0.73)	43.04 (-0.11)	38.61 (-4.54)	67.02 (+23.87)
	KM	11.95	13.28 (+1.33)	11.79 (-0.16)	11.78 (-0.17)	29.56 (+17.61)
	IFE	35.28	35.53 (+0.25)	35.03 (-0.25)	33.76 (-1.52)	61.87 (+26.59)
	RE_Math	94.99	95.71 (+0.72)	94.81 (-0.18)	95.17 (+0.18)	89.09 (-5.90)

Table 3: **Hallucination Rates Under Different Shot Settings.** All main values are in % (omitted). Deltas for each model are relative to its own baseline (1-shot). Only deltas are colored: **red** indicates improvement (lower), and **green** indicates degradation (higher). “-” denotes unavailable data.

Overall, common mitigation strategies are better understood as operations that impose bias or reinforcement on specific components. As a result, their benefits are strongly mechanism dependent and may come at the expense of other components. Reasoning SFT is more inclined to repair failures in reasoning integration, but it may amplify biases related to knowledge retrieval or constraint execution. Instruct models tend to improve the stability of instruction execution and alleviate some knowledge conflicts, but this does not equate to stronger multistep reasoning. Meanwhile, simply increasing or decreasing ICL demonstrations provides limited additional knowledge and contextual support, so the marginal effect remains small.

4.2 Causes of Hallucination Patterns (RQ3)

To further elucidate the mechanisms underlying model hallucinations in knowledge conflicts, we compare the Attention Maps of KE and KM in Figure 5.

Dominance of Parametric Priors. In KE cases, the input provides explicit evidence that contradicts the model’s parametric knowledge. The top of Figure 5 shows a diffuse attention pattern: the

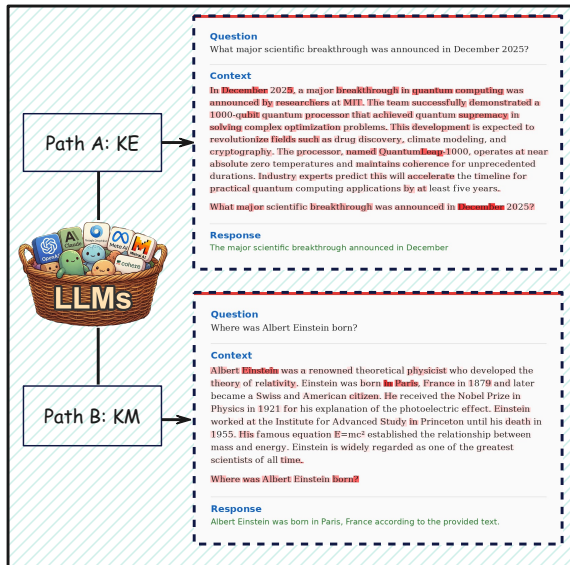


Figure 5: Visualization of Attention Maps in KE and KM

model attends to general entity tokens but does not strongly focus on the specific evidence tokens that determine the correct fact. As a result, strong parametric priors pull the model toward its memorized retrieval and suppress evidence extraction, causing the final answer to drift away from the provided facts.

Dominance of Misleading Context. KM cases fail differently. When the model lacks the needed parametric knowledge, or even when the knowledge exists but is not effectively retrieved, misleading context can dominate. The bottom of Figure 5 shows attention becoming overly concentrated on wrong tokens in the deceptive context. The Qwen3-4B example in Appendix I is consistent with this: instead of detecting an information gap or a conflict, the model treats misleading cues as evidence and produces confident but incorrect conclusions. This implies that without strong internal knowledge, the model is easily misled by false inputs, prioritizing consistency with the context over factual accuracy.

Computational Cost of Refusal. Table 4 quantifies the computational cost of getting this right. Across all four models, Qwen3-4B, Llama3.1-8B, Qwen3-1.7B, and Qwen3-14B, the refusal state in the IFE task consistently uses more ActN and higher GFLOPs than the non-refusal state. Refusal also shows lower attention sparsity, suggesting the model must attend to more signals and perform heavier cross-checking to identify conflicts and suppress hallucinations, whereas generating a hal-

Model	Task Type	ActN	AttnSp	GFLOPs
Qwen3-4B	IFE (Non-Refusal)	7.26 M	0.0191	215.7
	IFE (Refusal)	7.72 M	0.0175	229.3
LLama3.1-8B	IFE (Non-Refusal)	10.41 M	0.0148	468.9
	IFE (Refusal)	11.06 M	0.0137	498.4
Qwen3-1.7B	IFE (Non-Refusal)	4.73 M	0.0208	91.7
	IFE (Refusal)	5.03 M	0.0191	97.5
Qwen3-14B	IFE (Non-Refusal)	13.58 M	0.0169	755.0
	IFE (Refusal)	14.44 M	0.0154	802.6

Table 4: IFE Task Comparison on different answer, Abbreviations: ActN = active neurons; AttnSp = attention sparsity; GFLOPs = giga floating-point operations.

lucinated answer can be computationally cheaper.

5 Conclusion

We introduce *PRISM*, a benchmark designed to evaluate hallucination dimensions based on the response pipeline of LLMs. By decomposing the hallucination phenomenon into four dimensions known as KE, KM, RE, and IFE, we achieve identification of the stages where models fail. Extensive experiments demonstrate that while proprietary models generally outperform open-source models, all models exhibit significant trade-offs across different dimensions. Existing mitigation strategies function essentially as inductive biases imposed on specific generation stages. These approaches often improve performance in one dimension while compromising the stability of memory retrieval or logical reasoning. *PRISM* provides a quantitative basis for understanding these complex interactions and lays a solid foundation for future model selection, targeted training optimization, and trustworthy hallucination governance.

6 Limitations

Although *PRISM* provides a robust framework for hallucination diagnosis, this study still has several limitations. First, *PRISM* currently focuses solely on diagnosing hallucinations in the text modality and does not yet cover cross-modal hallucination challenges introduced by visual and auditory information in multimodal large language models. Meanwhile, the benchmark concentrates on five high-resource languages, and its applicability to low-resource linguistic settings remains to be validated. Second, real-world knowledge is continuously evolving, and static evaluation sets struggle to fully capture knowledge updating and forgetting

when models encounter real-time information. Although our OOD temporal knowledge data spans March–November 2025 and the training cutoffs of most evaluated models predate this window, we cannot entirely rule out the risk of data contamination for models whose training boundaries are undisclosed. Furthermore, the benchmark relies on authoritative static corpora to ensure verifiability and reproducibility, which may result in insufficient coverage of long-tail information.

Finally, *PRISM* is intended as a research diagnostic tool and should not serve as the sole basis for assessing model safety in high-stakes domains such as healthcare and law. Future work will extend evaluation to multimodal and multilingual settings, introduce rolling temporal window updates and systematic contamination detection, to further enhance the benchmark’s practicality and robustness.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (72442025). We thank the anonymous participants for taking part in our study. We are also grateful to the members of the DUFÉ Fintech Lab for their helpful comments.

References

- Samir Abdaljalil, Hasan Kurban, and Erchin Serpedin. 2025. [HalluVerse25: Fine-grained Multilingual Benchmark Dataset for LLM Hallucinations](#). *arXiv preprint arXiv:2503.07833*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2024. [GPT-4 Technical Report](#). *arXiv preprint arXiv:2303.08774*.
- AIMeta. 2025. [The Llama4 Herd: The Beginning of a New Era of Natively Multimodal AI Innovation](#).
- Aisha Alansari and Hamzah Luqman. 2026. [Large Language Models Hallucination: A Comprehensive Survey](#). *arXiv preprint arXiv:2510.06265*.
- Dang Anh-Hoang, Vu Tran, and Le-Minh Nguyen. 2025. [Survey and analysis of hallucinations in large language models: attribution to prompting strategies or model behavior](#). *Frontiers in Artificial Intelligence*, Volume 8 - 2025.
- Anthropic. 2025. [Claude3.7 Sonnet and Claude Code](#).
- Orlando Ayala and Patrice Bechard. 2024. [Reducing hallucination in structured outputs via Retrieval-Augmented Generation](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry Track)*, pages 228–238.
- Robert Baldock, Hartmut Maennel, and Behnam Neyshabur. 2021. [Deep Learning Through the Lens of Example Difficulty](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 10876–10889.
- Deborah L Bandalos. 2018. *Measurement theory and applications for the social sciences*. Guilford Publications.
- Yejin Bang, Ziwei Ji, Alan Schelten, Anthony Hartshorn, Tara Fowler, Cheng Zhang, Nicola Cancedda, and Pascale Fung. 2025. [HalluLens: LLM Hallucination Benchmark](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 24128–24156.
- Forrest Sheng Bao, Miaoran Li, Renyi Qu, Ge Luo, Erana Wan, Yujia Tang, Weisi Fan, Manveer Singh Tamber, Suleman Kazi, Vivek Sourabh, Mike Qi, Ruixuan Tu, Chenyu Xu, Matthew Gonzales, Ofer Mendelevitch, and Amin Ahmad. 2025. [FaithBench: A Diverse Hallucination Benchmark for Summarization by Modern LLMs](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 448–461.
- Adrien Barbaresi. 2021. [Trafalatura: A Web Scraping Library and Command-Line Tool for Text Discovery and Extraction](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 122–131.
- Bernd Bohnet, Vinh Q. Tran, Pat Verga, Roei Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, Tom Kwiatkowski, Ji Ma, Jianmo Ni, Lierni Sestorain Saralegui, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan Das, and 3 others. 2023. [Attributed Question Answering: Evaluation and Modeling for Attributed Large Language Models](#). *arXiv preprint arXiv:2212.08037*.
- Daniil A Boiko, Robert MacKnight, Gabe Gomes, and Ben Kline. 2023. [Autonomous chemical research with large language models](#). *Nature*, 624(7992):570–578.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldasari, Andrew D White, and Philippe Schwaller. 2024. [Augmenting large language models with chemistry tools](#). *Nature Machine Intelligence*, 6(5):525–535.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. [Language Models are Few-Shot Learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. 2005. Learning to Rank using Gradient Descent. In *Proceedings of the 22nd international conference on Machine learning*, pages 89–96.
- Neeloy Chakraborty, Melkior Ornik, and Katherine Driggs-Campbell. 2025. [Hallucination Detection in Foundation Models for Decision-Making: A Flexible Definition and Review of the State of the Art](#). *ACM Comput. Surv.*, 57(7).
- Kedi Chen, Qin Chen, Jie Zhou, He Yishen, and Liang He. 2024a. [DiaHalu: A Dialogue-level Hallucination Evaluation Benchmark for Large Language Models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9057–9079.
- Shiqi Chen, Yiran Zhao, Jinghan Zhang, I-Chun Chern, Siyang Gao, Pengfei Liu, and Junxian He. 2023. [FELM: benchmarking factuality evaluation of large language models](#). In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pages 44502–44523.
- Xiang Chen, Duanzheng Song, Honghao Gui, Chenxi Wang, Ningyu Zhang, Yong Jiang, Fei Huang, Chengfei Lyu, Dan Zhang, and Huajun Chen. 2024b. [FactCHD: Benchmarking Fact-Conflicting Hallucination Detection](#). In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 6216–6224.
- Qinyuan Cheng, Tianxiang Sun, Wenwei Zhang, Siyin Wang, Xiangyang Liu, Mozhi Zhang, Junliang He, Mianqiu Huang, Zhangyue Yin, Kai Chen, and Xipeng Qiu. 2023. [Evaluating Hallucinations in Chinese Large Language Models](#). *arXiv preprint arXiv:2310.03368*.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. [Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities](#). *arXiv preprint arXiv:2507.06261*.
- Jiaxi Cui, Munan Ning, Zongjian Li, Bohua Chen, Yang Yan, Hao Li, Bin Ling, Yonghong Tian, and Li Yuan. 2024. [Chatlaw: A Multi-Agent Collaborative Legal Assistant with Knowledge Graph Enhanced Mixture-of-Experts Large Language Model](#). *arXiv preprint arXiv:2306.16092*.
- Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. 2024. [Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models](#). *Journal of Legal Analysis*, 16(1):64–93.
- D. Diakoulaki, G. Mavrotas, and L. Papayannakis. 1995. [Determining objective weights in multiple criteria problems: The critic method](#). *Computers & Operations Research*, 22(7):763–770.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang (Lorraine) Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena Hwang, Soumya Sanyal, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. 2023. [Faith and Fate: Limits of Transformers on Compositionality](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 70293–70332.
- Deanna Emery, Michael Goitia, Freddie Vargus, and Iulia Neagu. 2025. [HalluMix: A Task-Agnostic, Multi-Domain Benchmark for Real-World Hallucination Detection](#). *arXiv preprint arXiv:2505.00506*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, and 1 others. 2020. [The Pile: An 800GB Dataset of Diverse Text for Language Modeling](#). *arXiv preprint arXiv:2101.00027*.
- Google DeepMind. 2026. [Gemini 3.1 pro - model card](#).
- Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Aditya K, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory Dickinson, Haggai Porat, Jason Hegland, and 21 others. 2023. [LegalBench: A Collaboratively Built Benchmark for Measuring Legal Reasoning in Large Language Models](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 44123–44279.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. [DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning](#). *Nature*, 645(8081):633–638.
- Stephen N Haynes, David Richard, and Edward S Kubany. 1995. [Content Validity in Psychological Assessment: A Functional Approach to Concepts and Methods](#). *Psychological assessment*, 7(3):238.
- Yancheng He, Shilong Li, Jiaheng Liu, Yingshui Tan, Weixun Wang, Hui Huang, Xingyuan Bu, Hangyu Guo, Chengwei Hu, Boren Zheng, Zhuoran Lin, Dekai Sun, Zhicheng Zheng, Wenbo Su, and Bo Zheng. 2025. [Chinese SimpleQA: A Chinese Factuality Evaluation for Large Language Models](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 19182–19208.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring Massive Multitask Language Understanding](#). In *9th International Conference on*

- Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.*
- Alex Heyman and Joel Zylberberg. 2025. [Reasoning Large Language Model Errors Arise from Hallucinating Critical Problem Features](#). *arXiv preprint arXiv:2505.12151*.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The Curious Case of Neural Text Degeneration](#). In *International Conference on Learning Representations*.
- Xiangkun Hu, Dongyu Ru, Lin Qiu, Qipeng Guo, Tianhang Zhang, Yang Xu, Yun Luo, Pengfei Liu, Yue Zhang, and Zheng Zhang. 2024. [RefChecker: Reference-based Fine-grained Hallucination Checker and Benchmark for Large Language Models](#). *arXiv preprint arXiv:2405.14486*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. [A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions](#). *ACM Transactions on Information Systems*, 43(2):1–55.
- Ziwei Ji, Yuzhe Gu, Wenwei Zhang, Chengqi Lyu, Dahua Lin, and Kai Chen. 2024. [ANAH: Analytical Annotation of Hallucinations in Large Language Models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8135–8158.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023a. [Survey of Hallucination in Natural Language Generation](#). *ACM Computing Surveys*, 55(12):1–38.
- Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023b. [Towards Mitigating LLM Hallucination via Self Reflection](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1827–1843.
- Nan Jiang, Qi Li, Lin Tan, and Tianyi Zhang. 2024a. [Collu-Bench: A Benchmark for Predicting Language Model Hallucinations in Code](#). *arXiv preprint arXiv:2410.09997*.
- Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang, Qun Liu, and Wei Wang. 2024b. [FollowBench: A Multi-level Fine-grained Constraints Following Benchmark for Large Language Models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4667–4688.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. [Large Language Models Struggle to Learn Long-Tail Knowledge](#). In *International conference on machine learning*, pages 15696–15707. PMLR.
- Michael Kane. 2013. [The Argument-Based Approach to Validation](#). *School Psychology Review*, 42(4):448–457.
- Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, and Kentaro Inui. 2023. [RealTime QA: What’s the Answer Right Now?](#) In *Advances in Neural Information Processing Systems*, volume 36, pages 49025–49043.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. [Decomposed Prompting: A Modular Approach for Solving Complex Tasks](#). In *The Eleventh International Conference on Learning Representations*.
- Emre Kiciman, Robert Ness, Amit Sharma, and Chenhao Tan. 2024. [Causal Reasoning and Large Language Models: Opening a New Frontier for Causality](#). *Transactions on Machine Learning Research*.
- Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmongkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, and 33 others. 2022. [Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets](#). *Transactions of the Association for Computational Linguistics*, 10:50–72.
- Yongchan Kwon, Shang Zhu, Federico Bianchi, Kaitlyn Zhou, and James Zou. 2025. [ReasonIF: Large Reasoning Models Fail to Follow Instructions During Reasoning](#). *arXiv preprint arXiv:2510.15211*.
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, and 1 others. 2023. [Measuring Faithfulness in Chain-of-Thought Reasoning](#). *arXiv preprint arXiv:2307.13702*.
- Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, and 35 others. 2023. [The BigScience ROOTS Corpus: A 1.6TB Composite Multilingual Dataset](#). *arXiv preprint arXiv:2303.03915*.
- Ernests Lavrinovics, Russa Biswas, Katja Hose, and Johannes Bjerva. 2025. [MultiHal: Multilingual Dataset for Knowledge-Graph Grounded Evaluation of LLM Hallucinations](#). *arXiv preprint arXiv:2505.14101*.
- Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023. [HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large](#)

- Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6449–6464.
- Yihan Li, Xiyuan Fu, Ghanshyam Verma, Paul Buitelaar, and Mingming Liu. 2025. **Mitigating Hallucination in Large Language Models (LLMs): An Application-Oriented Survey on RAG, Reasoning, and Agentic Systems.** *arXiv preprint arXiv:2510.24476*.
- Xun Liang, Shichao Song, Simin Niu, Zhiyu Li, Feiyu Xiong, Bo Tang, Yezhaohui Wang, Dawei He, Cheng Peng, Zhonghao Wang, and Haiying Deng. 2024a. **UHGEval: Benchmarking the Hallucination of Chinese Large Language Models via Unconstrained Generation.** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5266–5293.
- Xun Liang, Shichao Song, Zifan Zheng, Hanyu Wang, Qingchen Yu, Xunkai Li, Rong-Hua Li, Yi Wang, Zhonghao Wang, Feiyu Xiong, and Zhiyu Li. 2024b. **Internal Consistency and Self-Feedback in Large Language Models: A Survey.** *arXiv preprint arXiv:2407.14507*.
- Qian Lin, Junyi Li, and Hwee Tou Ng. 2025. **DynaQuest: A Dynamic Question Answering Dataset Reflecting Real-World Knowledge Updates.** In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 26918–26936.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. **TruthfulQA: Measuring How Models Mimic Human Falsehoods.** In *Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: long papers)*, pages 3214–3252.
- Yen-Ting Lin and Yun-Nung Chen. 2023. **LLM-EVAL: Unified Multi-Dimensional Automatic Evaluation for Open-Domain Conversations with Large Language Models.** In *Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)*.
- MingShan Liu and Jialing Fang. 2025. **Enhancing Mathematical Reasoning in Large Language Models with Self-Consistency-Based Hallucination Detection.** *arXiv preprint arXiv:2504.09440*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. **Lost in the Middle: How Language Models Use Long Contexts.** *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Ye Liu, Kazuma Hashimoto, Yingbo Zhou, Semih Yavuz, Caiming Xiong, and Philip S. Yu. 2021. **Dense Hierarchical Retrieval for Open-Domain Question Answering.** In *Findings of the Association for Computational Linguistics: EMNLP 2021*.
- Zhuang Liu, Degen Huang, Kaiyu Huang, Zhuang Li, and Jun Zhao. 2020. **FinBERT: A Pre-trained Financial Language Representation Model for Financial Text Mining.** In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 4513–4519. ijcai.org.
- Zhuang Liu, Shiyao Qian, Shuirong Cao, and Tianyu Shi. 2025. **Mitigating Age-related Bias in Large Language Models: Strategies for Responsible Artificial Intelligence Development.** *INFORMS Journal on Computing*, 2(1):1–22.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. **When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories.** In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nat McAleese. 2022. **Teaching language models to support answers with verified quotes.** *arXiv preprint arXiv:2203.11147*.
- Samuel Messick. 1994. **VALIDITY OF PSYCHOLOGICAL ASSESSMENT: VALIDATION OF INFERENCES FROM PERSONS’ RESPONSES AND PERFORMANCES AS SCIENTIFIC INQUIRY INTO SCORE MEANING.** *ETS Research Report Series*, 1994(2):i–28.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. **FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation.** In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100.
- Abhika Mishra, Akari Asai, Vidhisha Balachandran, Yizhong Wang, Graham Neubig, Yulia Tsvetkov, and Hannaneh Hajishirzi. 2024. **Fine-grained Hallucination Detection and Editing for Language Models.** *arXiv preprint arXiv:2401.06855*.
- Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, KaShun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. **RAGTruth: A Hallucination Corpus for Developing Trustworthy Retrieval-Augmented Language Models.** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10862–10878.
- Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, Zhenjiang Jin, Guang Liang, Rui Zhang, Wenzheng Zhang, Yuan Qu, Zhifei Ren, Yuefeng Sun, Yuanhong Zheng, Dongsheng Ma, and 42 others. 2025. **MinerU2.5: A Decoupled Vision-Language Model for Efficient High-Resolution Document Parsing.** *arXiv preprint arXiv:2509.22186*.
- Jio Oh, Soyeon Kim, Junseok Seo, Jindong Wang, Ruo Chen Xu, Xing Xie, and Steven Euijong Whang. 2024. **ERBench: An Entity-Relationship based Automatically Verifiable Hallucination Benchmark for**

- Large Language Models.** In *Advances in Neural Information Processing Systems*, volume 37, pages 53064–53101.
- OpenAI. 2024. Hello gpt-4o.
- OpenAI. 2025. Responses API Reference. <https://platform.openai.com/docs/api-reference/responses/create>. Accessed 2025-12-26.
- Jie Ouyang, Tingyue Pan, Mingyue Cheng, Ruiran Yan, Yucong Luo, Jiaying Lin, and Qi Liu. 2025. **HoH: A Dynamic Benchmark for Evaluating the Impact of Outdated Information on Retrieval-Augmented Generation.** In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. **Training language models to follow instructions with human feedback.** In *Advances in neural information processing systems*, volume 35, pages 27730–27744.
- Lawrence A Palinkas, Sarah M Horwitz, Carla A Green, Jennifer P Wisdom, Naihua Duan, and Kimberly Hoagwood. 2015. **Purposeful Sampling for Qualitative Data Collection and Analysis in Mixed Method Implementation Research.** *Administration and policy in mental health and mental health services research*, 42(5):533–544.
- Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. **Unifying large language models and knowledge graphs: A roadmap.** *IEEE Transactions on Knowledge and Data Engineering*, 36(7):3580–3599.
- Shrey Pandit, Jiawei Xu, Junyuan Hong, Zhangyang Wang, Tianlong Chen, Kaidi Xu, and Ying Ding. 2025. **MedHallu: A Comprehensive Benchmark for Detecting Medical Hallucinations in Large Language Models.** In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2858–2873.
- Bibek Paudel, Alexander Lyzhov, Preetam Joshi, and Puneet Anand. 2025. **HalluciNot: Hallucination Detection Through Context and Common Knowledge Verification.** *arXiv preprint arXiv:2504.07069*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. **Instruction Tuning with GPT-4.** *arXiv preprint arXiv:2304.03277*.
- Benjamin Plaut, Nguyen X. Khanh, and Tu Trinh. 2025. **Probabilities of Chat LLMs Are Miscalibrated but Still Predict Correctness on Multiple-Choice Q&A.** *Transactions on Machine Learning Research*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. **Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer.** *J. Mach. Learn. Res.*, 21(1).
- A B M Ashikur Rahman, Saeed Anwar, Muhammad Usman, and Ajmal Mian. 2024. **DefAn: Definitive Answer Dataset for LLMs Hallucination Evaluation.** *arXiv preprint arXiv:2406.09155*.
- Abhilasha Ravichander, Shruti Ghela, David Wadden, and Yejin Choi. 2025. **HALoGEN: Fantastic LLM Hallucinations and Where to Find Them.** In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1402–1425.
- Pranab Sahoo, Prabhaskar Meharia, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. 2024. **A Comprehensive Survey of Hallucination in Large Language, Image, Video and Audio Foundation Models.** In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11709–11724.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Abubakr Babiker, Nathanael Schärli, Aakanksha Chowdhery, Philip Mansfield, Dina Demner-Fushman, and 13 others. 2023. **Large language models encode clinical knowledge.** *Nature*, 620(7972):172–180.
- Samuel J Stratton. 2024. **Purposeful Sampling: Advantages and Pitfalls.** *Prehospital and Disaster Medicine*, 39(2):121–122.
- Atsushi Suzuki, Yulan He, Feng Tian, and Zhongyuan Wang. 2025. **Hallucinations are inevitable but can be made statistically negligible. the "innate" inevitability of hallucinations cannot explain practical llm issues.** *arXiv preprint arXiv:2502.12187*.
- Liyan Tang, Igor Shalyminov, Amy Wong, Jon Burnsky, Jake Vincent, Yu'an Yang, Siffi Singh, Song Feng, Hwanjun Song, Hang Su, Lijia Sun, Yi Zhang, Saab Mansour, and Kathleen McKeown. 2024. **TofuEval: Evaluating Hallucinations of LLMs on Topic-Focused Dialogue Summarization.** 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 4455–4480.
- GLM Team, Aohan Zeng, Bin Xu, and 1 others. 2024. **ChatGLM: A Family of Large Language Models from GLM-130B to GLM-4 All Tools.** *arXiv preprint arXiv:2406.12793*.
- Qwen Team. 2024. **QWQ: Reflect Deeply on the Boundaries of the Unknown.**
- Arun James Thirunavukkarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan,

- and Daniel Shu Wei Ting. 2023. [Large language models in medicine](#). *Nature Medicine*, 29(8):1930–1940.
- Konstantinos Thomas, Giorgos Filandrianos, Maria Lymperaiou, Chrysoula Zerva, and Giorgos Stamou. 2024. "i never said that": A dataset, taxonomy and baselines on response clarity classification. In *Findings of the Association for Computational Linguistics: EMNLP 2024*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: a Large-scale Dataset for Fact Extraction and VERification](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819.
- Yuchen Tian, Weixiang Yan, Qian Yang, Xuandong Zhao, Qian Chen, Wen Wang, Ziyang Luo, Lei Ma, and Dawn Song. 2025. [CodeHalu: Investigating Code Hallucinations in LLMs via Execution-based Verification](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 25300–25308.
- SMTI Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. [A Comprehensive Survey of Hallucination Mitigation Techniques in Large Language Models](#). *arXiv preprint arXiv:2401.01313*.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. [Language Models Don't Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 74952–74965.
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2024. [Fresh-LLMs: Refreshing Large Language Models with Search Engine Augmentation](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13697–13720.
- Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen-tse Huang, Pinjia He, Wenxiang Jiao, and Michael Lyu. 2024. [LogicAsker: Evaluating and Improving the Logical Reasoning Ability of Large Language Models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2124–2155.
- Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, Bo Zhang, Liqun Wei, Zhihao Sui, Wei Li, Botian Shi, Yu Qiao, Dahua Lin, and Conghui He. 2024a. [MinerU: An Open-Source Solution for Precise Document Content Extraction](#). *arXiv preprint arXiv:2409.18839*.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. 2024b. [A survey on large language model based autonomous agents](#). *Frontiers of Computer Science*, 18(6):186345.
- Jakob Warkotsch. 2018. Developing a graphical user interface for generating wikipedia lists with wikidata.
- Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. 2024a. [Measuring short-form factuality in large language models](#). *arXiv preprint arXiv:2411.04368*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-Thought Prompting Elicits Reasoning in Large Language Models](#). *arXiv preprint arXiv:2201.11903*.
- Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Jie Huang, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, Cosmo Du, and Quoc V. Le. 2024b. [Long-form factuality in large language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 80756–80827.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2020. [CCNet: Extracting High Quality Monolingual Datasets from Web Crawl Data](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4003–4012.
- Jiarui Wu, Zhuo Liu, and Hangfeng He. 2025a. [Mitigating Hallucinations in Multimodal Spatial Relations through Constraint-Aware Prompting](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 3450–3468.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2023. [AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation](#). *arXiv preprint arXiv:2308.08155*.
- Yuhe Wu, Yuran Chen, Zhuang Liu, and Wayne Lin. 2025b. [Enhancing Financial Decision-making under Cyber Threats: a Dual-branch Framework Integrating Bayesian Deep Learning and Explainable AI](#). *Annals of Operations Research*, pages 1–33.
- xAI. 2025. [Models](#). Accessed:2025-07-09.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, and 9 others. 2025. [The rise and potential of large language model based agents: A survey](#). *Science China Information Sciences*.

- Yi Xu, Junjie Ou, Hui Xu, and Luoyi Fu. 2023. [Temporal Knowledge Graph Reasoning with Historical Contrastive Learning](#). In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 4765–4773.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498.
- Jianhao Yan, Yun Luo, and Yue Zhang. 2024. [RefuteBench: Evaluating Refuting Instruction-Following for Large Language Models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13775–13791.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380.
- Lei Yu, Meng Cao, Jackie CK Cheung, and Yue Dong. 2024. [Mechanistic Understanding and Mitigation of Language Model Non-Factual Hallucinations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7943–7956.
- Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023. [Investigating the Catastrophic Forgetting in Multimodal Large Language Models](#). *arXiv preprint arXiv:2309.10313*.
- Cong Zhang, Xin Deik Goh, Dexun Li, Hao Zhang, and Yong Liu. 2025a. [Planning with Multi-Constraints via Collaborative Language Agents](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10054–10082.
- Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. 2024a. [Chain of Preference Optimization: Improving Chain-of-Thought Reasoning in LLMs](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 333–356.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2025b. [Siren’s Song in the AI Ocean: A Survey on Hallucination in Large Language Models](#). *Computational Linguistics*, 51(4):1373–1418.
- Yuji Zhang, Sha Li, Jiateng Liu, Pengfei Yu, Yi R Fung, Jing Li, Manling Li, and Heng Ji. 2024b. [Knowledge Overshadowing Causes Amalgamated Hallucination in Large Language Models](#). *arXiv preprint arXiv:2407.08039*.
- Yuxiang Zhang, Jing Chen, Junjie Wang, Yaxin Liu, Cheng Yang, Chufan Shi, Xinyu Zhu, Zihao Lin, Hanwen Wan, Yujiu Yang, Tetsuya Sakai, Tian Feng, and Hayato Yamana. 2024c. [ToolBeHonest: A Multi-level Hallucination Diagnostic Benchmark for Tool-Augmented Large Language Models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11388–11422.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. [Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623.
- Zhiying Zhu, Yiming Yang, and Zhiqing Sun. 2024. [HaluEval-Wild: Evaluating Hallucinations of Language Models in the Wild](#). *arXiv preprint arXiv:2403.04307*.

Appendix

A Benchmark Comparisons

Table 5 provides a comprehensive comparison between *PRISM* and existing benchmarks in *Evaluation Scope* and *Methodological Design*. Most existing benchmarks concentrate on only certain aspects of hallucination, failing to provide a holistic examination of failures. In terms of the depth of the evaluation, they rely largely on posterior analysis, making the attribution of hallucination causes unclear.

In contrast, *PRISM* introduces a diagnostic framework that covers four distinct failure mechanisms. By defining hallucination categories upfront at the input level, *PRISM* enables a more accurate identification of where failures occur and why they arise.

B Related Works

B.1 Evolution of Hallucination Categorization

Early research typically classified hallucinations into general categories, such as Intrinsic versus Extrinsic types based on their relationship to the input (Ji et al., 2023a), or distinguished between Factuality and Faithfulness errors based on consistency with external knowledge (Zhang et al., 2025b; Huang et al., 2025). With the expansion of model capabilities in 2024, recent studies have extended these definitions to specific domains. For instance, (Dahl et al., 2024) analyzed premise-compliance errors in legal tasks, where models erroneously validate incorrect user assumptions. Furthermore, theoretical studies (Sahoo et al., 2024; Suzuki et al., 2025) suggest that hallucinations are inherent statistical features of probabilistic models rather than simple engineering defects.

However, existing studies mostly focus on the manifestation of errors rather than the root causes. Current frameworks often fail to distinguish whether a failure stems from missing data, flawed reasoning, or an inability to follow instructions, thereby hindering the accurate attribution of errors in complex generation pipelines.

B.2 Benchmarks for Hallucination Evaluation

In response to the challenge of hallucination in LLMs, previous studies have developed a series of benchmarks to evaluate hallucinations. Early efforts primarily assessed the model’s grasp of static world knowledge and its resilience to interference.

Specifically targeting hallucinations, HaluEval (Li et al., 2023) generates and screens a large number of hallucination samples to evaluate whether LLMs can identify hallucination issues. As understanding of hallucination mechanisms deepened, the evaluation paradigms shifted toward higher granularity. FActScore (Min et al., 2023) introduced an atomic-level evaluation, dividing long-form text into individual facts to verify their support against reliable knowledge sources. Similarly, FELM (Chen et al., 2023) adopted segment-level annotation, significantly improving the precision of error localization across diverse domains. Recent frameworks have further systematized the definition of errors: HaluLens (Bang et al., 2025) formalized the distinction between extrinsic hallucinations (contradicting training data or reality) and intrinsic hallucinations (deviating from input context), while FollowBench (Jiang et al., 2024b) and HALoGEN (Ravichander et al., 2025) expanded the evaluation scope to include instruction following failures and source-based error attribution, distinguishing memory distortion from fabrication.

Despite these advances, as shown in Table 5, existing benchmarks suffer from a fundamental methodological limitation: their reliance on outcome-oriented and posterior evaluation. *PRISM* is designed to address this critical gap. Unlike prior works, *PRISM* establishes a diagnostic framework based on active probing, isolating the distinct cognitive stages of source memory retrieval, reasoning, and instruction following to precisely locate the origin of hallucinations within the generative pipeline.

B.3 Strategies for Hallucination Mitigation

Numerous studies have proposed diverse approaches to mitigate hallucinations in LLMs (Tonmoy et al., 2024). From the model generation perspective, existing hallucination mitigation methods target hallucination mechanisms at different stages, primarily memory, reasoning, and instruction-following (Li et al., 2025). Memory-stage hallucinations are commonly addressed through knowledge-centric approaches, including retrieval-augmented generation (Ayala and Bechard, 2024) and knowledge base querying (Pan et al., 2024), which aim to alleviate errors caused by missing or incorrect knowledge. Reasoning-stage hallucinations, stemming from inconsistencies in multi-step reasoning, are mitigated by structuring the reasoning process via techniques such as Chain-of-Thought prompting (Zhang et al., 2024a) and

Benchmark	Evaluation Scope				Methodological Design	
	KE	KM	RE	IFE	Variable Control	Diag. Mode
TruthfulQA (Lin et al., 2022)	✓	✗	✗	✗	✗	✗
HaluEval (Li et al., 2023)	✓	✓	✓	✗	✗	✗
FactScore (Min et al., 2023)	✓	✓	✗	✗	✗	✗
FactCHD (Chen et al., 2024b)	✓	✗	✓	✗	✗	✗
FELM (Chen et al., 2023)	✓	✗	✓	✗	✗	✗
HalluQA (Cheng et al., 2023)	✓	✗	✗	✓	✗	✓
FreshQA (Vu et al., 2024)	✓	✓	✗	✗	✓	✗
FollowBench (Jiang et al., 2024b)	✗	✗	✓	✓	✗	✗
SimpleQA (Wei et al., 2024a)	✓	✓	✗	✗	✗	✗
Collu-Bench (Jiang et al., 2024a)	✗	✗	✓	✓	✓	✗
LongFact (Wei et al., 2024b)	✓	✗	✗	✗	✗	✗
RAGTruth (Niu et al., 2024)	✓	✗	✓	✗	✗	✗
DiaHalu (Chen et al., 2024a)	✗	✗	✓	✓	✗	✗
LogicAsker (Wan et al., 2024)	✗	✗	✓	✗	✓	✓
RefuteBench (Yan et al., 2024)	✗	✗	✓	✓	✗	✗
HaluEval-Wild (Zhu et al., 2024)	✓	✓	✓	✗	✗	✓
REFCHECKER (Hu et al., 2024)	✗	✗	✓	✓	✗	✗
ERBench (Oh et al., 2024)	✗	✗	✓	✓	✗	✗
ToolBeHonest (Zhang et al., 2024c)	✗	✗	✓	✓	✗	✗
Defan (Rahman et al., 2024)	✓	✗	✓	✓	✗	✗
UHGEval (Liang et al., 2024a)	✓	✗	✓	✗	✗	✗
WildHallucinations (Liang et al., 2024a)	✗	✓	✓	✗	✗	✗
TofuEval (Tang et al., 2024)	✓	✓	✓	✗	✗	✗
FavaBench (Mishra et al., 2024)	✓	✓	✓	✗	✗	✗
ANAH (Ji et al., 2024)	✓	✓	✓	✗	✗	✗
Chinese SimpleQA (He et al., 2025)	✓	✓	✗	✗	✗	✗
HalluMix (Emery et al., 2025)	✗	✓	✓	✗	✗	✗
MedHallu (Pandit et al., 2025)	✓	✓	✗	✗	✗	✗
HalluVerse25 (Abdaljalil et al., 2025)	✓	✗	✓	✓	✗	✗
ReasonIF (Kwon et al., 2025)	✗	✗	✓	✓	✗	✓
MultiHal (Lavrinovics et al., 2025)	✓	✗	✓	✗	✗	✓
DynaQuest (Lin et al., 2025)	✓	✓	✗	✗	✓	✗
CodeHalu (Tian et al., 2025)	✓	✗	✓	✗	✗	✓
FaithBench (Bao et al., 2025)	✗	✗	✓	✓	✗	✗
HalluciNot (Paudel et al., 2025)	✓	✗	✓	✗	✗	✗
HALoGEN (Ravichander et al., 2025)	✓	✓	✓	✗	✗	✗
HalluLens (Bang et al., 2025)	✓	✓	✓	✗	✗	✗
PRISM (Ours)	✓	✓	✓	✓	✓	✓

Table 5: Systematic Comparison of State-of-the-Art Hallucination and Trustworthiness Benchmarks (2022–2025).

Note: **Variable Control** denotes whether a benchmark enforces orthogonal isolation of causal factors, ensuring attribution correctness by preventing confounding effects when decomposing hallucination sources.

self-consistency reasoning (Liu and Fang, 2025), often combined with causal learning (Kiciman et al., 2024) or contrastive learning (Xu et al., 2023) to improve reasoning robustness. Instruction-stage hallucinations, characterized by deviations from user instructions or task constraints, are typically mitigated through constrained prompting (Wu et al., 2025a) and post-hoc verification mechanisms, including self-reflection (Ji et al., 2023b) and self-consistency checking (Liang et al., 2024b), to enhance instruction-following behavior.

While prior work studies hallucinations from multiple sources, systematic evaluation across different hallucination problems is limited (Anh-Hoang et al., 2025; Chakraborty et al., 2025).

Therefore we introduce **PRISM** and use it to analyze mitigation strategies across hallucinations caused by memory, reasoning, and instruction-following failures.

C Data Source and Task Definitions

In this section, we present the comprehensive taxonomy of the **PRISM** benchmark, detailing the specific failure mechanisms associated with each cognitive stage. To facilitate a granular diagnosis of model hallucinations, we categorize failures into four primary dimensions: KE, KM, RE, and IFE. Each dimension is further divided into specific sub-categories to capture distinct error patterns ranging from factual distortions to procedural breakdowns.

Table 6 provides rigorous definitions for these sub-categories and maps them to their respective data sources, illustrating how our multi-source construction strategy ensures broad coverage across the entire spectrum of potential failures.

For the temporal OOD portion of KM, our TimeOOD data were collected between March and November 2025, which is strictly later than the officially disclosed knowledge cutoffs or public time-boundary information of several representative models evaluated, including Llama 4 Scout (AIMeta, 2025), the Grok 4 series (xAI, 2025), and Gemini 3 Pro (Google DeepMind, 2026).

D Construction Pipeline

To ensure high data quality, we implement a rigorous four-stage construction pipeline. As summarized in Table 7, this process filters an initial pool of over 33,000 candidate samples down to the final 9,448 high-quality instances in the PRISM benchmark. Full implementation details for each stage are provided in the following subsections.

D.1 Data Preprocessing

Data Conversion. Our self-constructed raw materials include both web-based text and PDF documents. To unify downstream processing, we convert materials from different sources into a sequential Markdown representation. For web-based text, we use Trafilatara (Barbaresi, 2021) for main-content extraction; this tool is designed for boilerplate removal in web documents and effectively separates core content from page noise. For PDF files, we use MinerU (Wang et al., 2024a; Niu et al., 2025). This tool understands complex page layouts, allowing it to accurately reconstruct the document’s structure and maintain the correct reading flow.

Data Cleaning. After document conversion, we apply rule-based data cleaning to remove noise that is not directly relevant to evaluation instance construction. First, we eliminate overlapping or highly similar text segments to reduce redundancy and prevent repeated information from affecting downstream processing. Second, we remove in-text citation markers, footnotes, and external links, and discard descriptive or metadata content such as copyright notices, navigation cues, and residual formatting artifacts. Finally, text segments that are clearly incomplete, lack sufficient context, or cannot be understood independently are removed

at this stage. Based on these steps, we obtain a consistent and readable raw corpus.

Corpus Formation. We construct pre-defined text units as the basic processing granularity for downstream multi-agent workflows. Prior work shows that full-document processing leads to overly long contexts and poor localization, while paragraph-level inputs often lack sufficient context; semantically coherent text units offer a better trade-off between contextual completeness and efficiency (Liu et al., 2021, 2024). Accordingly, after document conversion and data cleaning, text units are built from section hierarchy and paragraph boundaries, merging adjacent paragraphs when necessary and applying token-level length constraints. For web documents without explicit structure, units are formed using paragraph boundaries. Each unit preserves source identifiers and positional metadata to enable precise localization and traceability.

D.2 Multi-agent Data Construction

Generating complex evaluation samples in a single pass is often unreliable because models struggle to follow multiple rules at the same time (Wei et al., 2023). To address this, we propose a four-agent framework that breaks the construction process into smaller steps. Unlike black box generation, our method transforms raw text into structured data stage by stage (Khot et al., 2023). We assign specific roles to separate agents (Wu et al., 2023). This ensures that each instance includes a clear question, valid evidence from the source, and a precise error label. This modular design automates the work and makes it easier to fix or improve specific parts of the process.

Schema Normalizer Agent. The Schema Normalizer Agent converts the raw corpus into a unified standard format. It separates the question, answer, and source information, without adding new facts. The standardized output then serves as a basis for the Evidence Retriever Agent, allowing it to focus on finding relevant supporting materials. The prompt used for this agent is shown in Appendix K.2 Prompt for Schema Normalizer.

Evidence Retriever Agent. The Evidence Retriever Agent ensures that every data instance is supported by clear evidence. We use Gemini 3 Pro, known for its strong long-context capabilities, to power this agent. It searches the raw corpus for specific text segments that validate the connec-

Subcategory	Definition	Source Category				
		EB	ED	RWK	HE	SD
Knowledge Error(KE)						
KE1: Factual Distortion (FD)	The model generates factually incorrect outputs despite possessing relevant knowledge, due to errors in knowledge representation, updating or retrieval.	✓		✓		
KE2: Intra-Memory Conflict (IMC)	The model internally stores multiple conflicting versions of the same fact, resulting in inconsistent or contradictory answers across different queries, contexts, or interaction turns.	✓	✓			
KE3: Entity-Identity Confusion (EIC)	The model fails to correctly distinguish between entities with identical names, multiple meanings, or semantic proximity, incorrectly transferring or binding knowledge from one entity to another.	✓	✓	✓		
Knowledge Missing(KM)						
KM1: Domain-Specific Knowledge (DSK)	The model lacks newer or niche knowledge in certain specific domains, which prevents it from answering related questions.	✓	✓	✓		
KM2: Fictional Knowledge (FK)	The model’s dataset does not include any knowledge of the fictional content in the question, which prevents the model from answering the relevant question.			✓		✓
KM3: Timely Knowledge (TK)	The model’s dataset does not include content that changes continuously over time or recent events, which prevents the model from answering related questions.			✓		✓
KM4: Non-Public Knowledge (NPK)	When the model asks questions involving personal thoughts or private information, the dataset may not contain the question content or may refuse to answer due to privacy concerns, thus preventing the model from answering the relevant questions.			✓		✓
Reasoning Error(RE)						
RE1: Logical Fallacy (LF)	The model fails to adhere to abstract formal logic or causal rules, deriving invalid conclusions from valid premises.	✓	✓		✓	
RE2: Procedural Failure (PF)	The model fails to maintain the correct sequence or continuity in multi-step or multi-hop tasks, resulting in lost chains of thought or skipped intermediate steps.	✓	✓			
RE3: Information Integration Failure (IIF)	The model fails to identify, prioritize, or synthesize scattered information from long or noisy contexts, leading to reasoning based on incomplete evidence.	✓		✓		✓
RE4: Mathematical Reasoning Failure (MRF)	Errors in mathematical or algorithmic reasoning involving derivation, calculation, or implementation.	✓			✓	
Instruction Following Error(IFE)						
IFE1: Explicit Format (EF)	The model fails to adhere to strict structural specifications or text syntactic patterns, resulting in invalid data formats or template violations.	✓				
IFE2: Length Constraints (LC)	The model violates quantitative length requirements, failing to keep the output within the specified word count or boundary limits.	✓	✓			
IFE3: Language Constraints (LgC)	The model fails to respond in the specified target language, or incorrectly mixes multiple languages when a single language is required.	✓	✓			
IFE4: Complex & Cognitive Load (CCL)	The model fails to parse or execute complex instructions involving negation or conditional logic, often ignoring critical constraints or specific rule details.	✓			✓	

Table 6: Taxonomy of our benchmark and corresponding data sources. **EB**: Existing Benchmarks; **ED**: Enhanced Datasets derived from existing benchmarks; **RWK**: Real-world Knowledge collected from authoritative sources; **HE**: Human Exams serving as human-level reasoning references; **SD**: Synthetic Data generated under controlled constraints.

Stage	Action	KE	KM	RE	IFE	Total
Data Collection	Mining & Synthesis	4,982	8,117	10,696	9,539	33,334
Data Cleaning	Conversion & Denoising	4,200	6,824	9,075	7,943	28,042
Multi-agent Construction	Evidence Grounding & Quality Scoring	2,655	2,979	4,638	4,310	14,582
Human Selection	Expert Adjudication	1,933	2,078	2,995	2,442	9,448
Pass Rate	Final / Initial	38.80%	25.60%	28.00%	25.60%	28.34%

Table 7: Per-stage sample counts and pass rates of the PRISM construction pipeline

tion between the question and the answer. Any instance that relies on hidden assumptions or cannot be traced to a clear source in the text is removed at this stage. This ensures that the retained data is fully grounded, preventing the evaluation from relying on unverifiable background knowledge (Menick et al., 2022; Bohnet et al., 2023). The details of its prompt are in Appendix K.2 Prompt for Evidence Retriever.

Type Classifier Agent. The Type Classifier assigns each verified example to one of four failure dimensions. It reads the question and the reference answer, identifies what the example mainly tests, such as source memory, reasoning, or following instructions, and assigns one label. This step does not change the content. It only groups examples by dimension so we can report where models tend to fail. Appendix K.2 Prompt for Type Classifier presents the detailed prompt design.

Quality Scoring Agent. The Quality Scoring Agent evaluates each instance across multiple dimensions, focusing on factuality, discriminability, and clarity. Based on these criteria, the agent assigns a quality score to each instance. The primary goal of this step is to filter out low-quality data: only high-scoring instances are retained, while those that are trivial are discarded. This rigorous screening ensures that only challenging instances proceed to the human selection stage. The detailed prompt structure is presented in Appendix K.2 Prompt for Quality Scoring.

D.3 Human Selection

We propose an expert-driven, three-dimensional quality selection strategy for *PRISM* candidate samples. This strategy targets two task families formed by four dimensions in the benchmark, namely knowledge-related and reasoning/instruction-related tasks. The core goal is to ensure that the final item bank satisfies Clarity, Domain Relevance, and Coherence. The selection process proceeds in two stages: (i) strati-

fied purposive sampling and expert recruitment; (ii) multi-round review and filtering based on the three dimensions.

Expert Recruitment. We use stratified purposive sampling to form two non-overlapping expert review teams, and each team is confined to its corresponding task family to reduce the risk of cross-mechanism misjudgment and evaluation standard drift (Palinkas et al., 2015; Stratton, 2024). We recruit a total of $N = 80$ reviewers, including:

- **Panel A: Knowledge & Source Panel.** $n = 40$. Responsible for KE and KM samples. The core responsibility is to judge, without introducing external background knowledge, whether an item truly relies on parametric knowledge.
- **Panel B: Reasoning & Instruction Panel.** $n = 40$. Responsible for RE and IFE samples. The core responsibility is to ensure that any prerequisite knowledge involved in the sample is assumed to be known to the model, that the reasoning chain can be closed within the item itself, and that instruction constraints are clear and executable.

Candidates are sourced from universities and research institutes, NLP/IR/reasoning evaluation research communities, industry research and evaluation teams, and professional technical communities. Eligibility and screening criteria are set by panel: the Knowledge & Source Panel requires an advanced degree or three years or more of work experience in fact checking or information retrieval; the Reasoning & Instruction Panel requires three years or more of work experience in reasoning tasks or instruction specification development.

Filtering Rules. To implement our dual-track validation framework, we conduct a multi-round expert review process. Each candidate item is evaluated against a predefined set of track-specific quality dimensions. Experts provide independent scores

and brief rationales, and items are assigned one of three decisions: retain, remove, or revise then re-review. This iterative process continues until convergence, ensuring that every item in the final *PRISM* bank meets our quality baseline. The following sections detail the evaluation dimensions.

- **Clarity:** Clarity is the foundation of reliable selection. Once an item is written vaguely, a model’s score can be shifted by reading habits and by how it fills in missing premises, rather than by the ability differences we care about. Bandalos (Bandalos, 2018) points out that ambiguity introduced by item wording can introduce construct-irrelevant variance, thereby harming reliability and validity. In benchmark construction, clarity is also directly tied to face validity, namely whether readers can immediately see what the item is asking, whether the information is written explicitly, and whether answering does not require guessing the author’s intent. Only when clarity is met do the subsequent judgments of relevance and coherence become meaningful.

Core Question: *Is the item statement precise, or could readers arrive at two equally reasonable but different interpretations?*

- **Domain Relevance:** Domain relevance means that an item genuinely operates on the dimensions we define, rather than being driven by side issues. Haynes et al. (Haynes et al., 1995) summarize the key points of content validity in two aspects: the content should stay close to what is intended to be measured, and the coverage should be representative. In human selection, this means avoiding items that appear to test a certain dimension but end up mainly testing something else, for example, succeeding through obscure background knowledge, or creating differences through wording tricks. Items with strong domain relevance are less likely to be derailed by test-gaming patterns, and they better support interpreting errors as deficiencies in a specific dimension, which is consistent with the testing standards that emphasize evidence for score interpretation and use (Messick, 1994).

Core Question: *Does the main difficulty of this item primarily come from the dimension it is labeled with, rather than from irrelevant factors?*

- **Coherence:** Coherence requires internal consistency and a closed information chain. The information provided in the prompt should support the reference answer or label; the item should not present one story in the material and another in the answer, and it should not contain facts that conflict across sentences. Kane’s (Kane, 2013) argument-based approach to validation emphasizes that any interpretation relies on a chain of inferences; if the internal chain of an item breaks, it effectively pushes key inferences onto the model or the evaluator, making conclusions unstable. Similarly, FEVER binds support or refute judgments to necessary evidence so that the judgment does not drift away from the evidence (Thorne et al., 2018). We include coherence as a dimension to ensure that each sample can self-consistently explain why the answer is what it is and why the label is what it is.

Core Question: *Using only the information provided in the item, can one stably reach the same conclusion, and is the answer or label consistent with the material?*

E Question Quality Filter

E.1 Sample Quality Evaluation Dimension

We argue that a high-quality QA instance should exhibit factual correctness, unambiguous categorical targeting, and linguistic clarity. Accordingly, we define three scoring dimensions to evaluate sample quality: A (Factuality), B (Discriminability), and C (Clarity). More detailed criteria for each scoring dimension are provided in Table 8. Below, we outline the purpose and rationale of each dimension.

Factuality. Factuality evaluates whether an answer is correct and strictly grounded in the original source evidence. Classic benchmarks such as HotpotQA explicitly incorporate factuality by requiring models to answer multi-hop questions while identifying supporting evidence (Yang et al., 2018). Subsequent results on TruthfulQA show that even large-scale language models frequently produce imitative falsehoods, highlighting the importance of evaluating factual correctness (Lin et al., 2022). The scoring agent compares each claim in the answer against the original corpus to ensure evidence consistency, preventing the introduction of hallucinated content that is inconsistent with or unverifiable from the source material. A high score

A. Factuality	
1-2	Largely incorrect or contradictory to the source with no supporting source.
3-4	Core claims unsupported or inconsistent with the source with at most one weakly related source.
5-6	Some alignment with the source with main claims partially supported by one source.
7-8	Minor inaccuracies with main claims supported by two sources.
9-10	Fully faithful to the source with all claims supported by more than two sources.
B. Discriminability	
1-2	Absolutely unclassifiable.
3-4	Features of multiple hallucination categories present.
5-6	Dominant hallucination mechanism not immediately evident.
7-8	Largely attributable to a primary hallucination category.
9-10	No cross-category features.
C. Clarity	
1-2	Disorganized, vague, and difficult to interpret.
3-4	Noticeable ambiguity or redundancy.
5-6	Generally understandable.
7-8	Clear and fluent, with minor room for refinement.
9-10	Highly clear and precise, immediately interpretable.

Table 8: Evaluation criteria used in our study. Each dimension is rated on a 1–10 scale with detailed scoring guidelines.

indicates that the answer is factual, verifiable, and free from factual errors, logical contradictions, or unsupported assertions.

Discriminability. Discriminability measures whether the question precisely targets a single failure category, minimizing label overlap or semantic ambiguity. Baldock et al. (Baldock et al., 2021) introduced the SoftGap metric in OOD detection, showing that a larger margin implies more confident, unambiguous predictions. Plaut et al. (Plaut et al., 2025) further validate margin as a reliable uncertainty indicator across QA benchmarks, reinforcing its utility for measuring categorical precision. This margin-based approach ensures the rigor of our discriminability mechanism. Let s_i be the predicted confidence score for category i , computed via softmax:

$$s_i = \frac{\exp(z_i)}{\sum_k \exp(z_k)},$$

and let $s_{(1)}, s_{(2)}$ denote the highest and second-highest scores respectively. The margin is computed as:

$$\text{Margin} = \max_i s_i - \max_{j \neq i} s_j.$$

A larger margin reflects clear category targeting, while a smaller one indicates ambiguity. High discriminative precision ensures each QA pair isolates a specific failure mode, strengthening the benchmark’s effectiveness in identifying the root causes of hallucination.

Clarity. Clarity is the basis of valid assessment. We evaluate clarity for the question–answer pair as a whole. This means explicitly naming all entities and context and using precise, direct phrasing in well-structured sentences. When both the question and the answer are clearly phrased, there is effectively a single intended reading, so any factual inconsistency or hallucinated content can be unambiguously identified and traced. Indeed, discourse studies of evasive or ambiguous answers find that unclear responses often contain contradictions, topic shifts, or incomplete fragments, resulting in multiple interpretations (Liu et al., 2020; Thomas et al., 2024). This approach matches standard QA dataset practices. By enforcing clarity, we avoid these confounds and ensure that hallucination judgments reflect genuine content errors rather than mere linguistic confusion (Ouyang et al., 2025).

E.2 Sample Quality Score

We adopt a two-step filtering process to automatically assess and refine QA instances before final use. The first step is weighted evaluation, which aggregates factuality, discriminative precision and clarity into a weighted quality score. The second step applies threshold-based elimination to eliminate samples with significant flaws in any individual aspect. Both procedures are implemented by the Quality Scoring Agent to ensure scoring consistency, as detailed in the following paragraphs.

1-2 Very poor Factuality	Source Text
<p>Question: What was the global methane emission trend in 2019?</p> <p>Reference Answer: Global emissions increased in 2019 as satellite systems had not yet been launched and policies were not enforced.</p> <p>Assessment: Contradicts Section 3.2. Factual inaccuracy on both trend and timeline. Fabricated content regarding satellite systems.</p> <p>Supporting Evidence Count: 0 →Label: A (Factuality 1–2)</p>	<p>...Section 3.2: Global methane emissions declined by 11% between 2015 and 2024, attributed to updated regulatory frameworks introduced in North America and Oceania... The North American Methane Accord (NAMA) of 2020 introduced annual audit mechanisms, while Oceania’s regional climate pacts imposed sector-specific leak detection standards across energy infrastructure...</p> <p>...Section 4.1: ...Data quality assessments in South Asia flagged inconsistencies between satellite-based and ground-based methane monitoring. Until late 2021, satellite platforms lacked consistent cross-calibration protocols, introducing potential biases in longitudinal emission trends...</p> <p>...Section 4.3: The deployment of MethaneScan-X2 in Q1 2023 significantly enhanced spatial granularity, enabling real-time detection of localized emission clusters, particularly in previously under-monitored territories such as equatorial Africa and Central Asia...</p> <p>...Section 4.4: While the aggregate emissions trajectory remains downward as of mid-2025, substantial interregional variation persists. Modeling uncertainties stem from disparate meteorological datasets, platform-specific sensor drift, and incomplete archival baselines in low-income reporting zones...</p>
<p>3-4 Low Factuality</p> <p>Question: Summarize the main reasons for methane reduction from 2010 to 2020.</p> <p>Reference Answer: The drop was primarily due to changes in agriculture and global warming awareness across Asia and Africa.</p> <p>Assessment: Partially relevant. Fails to mention regulation (Section 3.2)</p> <p>Supporting Evidence Count: 1 →Label: A (Factuality 3–4)</p>	
<p>5-6 Medium Factuality</p> <p>Question: In what ways did MethaneScan-X contribute to the improvement of methane emissions tracking in previously underreported regions, and what specific geographical area benefited most from its deployment?</p> <p>Reference Answer: MethaneScan-X improved tracking precision in developing regions, particularly South Asia, by detecting smaller-scale emission sources.</p> <p>Assessment: Correct attribution to Section 4.1, but omits 2018 timestamp and link to Section 3.2’s downward trend.</p> <p>Supporting Evidence Count: 1–2 →Label: A (Factuality 5–6)</p>	
<p>7-8 High Factuality</p> <p>Question: How did the absence of consistent satellite calibration prior to 2016 affect methane emission measurements in South Asia, and what were the consequences for global methane datasets in that period?</p> <p>Reference Answer: Before 2016, missing satellite calibration in South Asia led to incomplete methane reporting, limiting the global dataset’s coverage. After 2016, improved calibration practices enabled more reliable regional estimates.</p> <p>Assessment: Could further connect regional effects to the emission decline discussed in Section 3.2, 4.1, 4.2, but overall precise.</p> <p>Supporting Evidence Count: 2–3 →Label: A (Factuality 7–8)</p>	
<p>9-10 High Factuality</p> <p>Question: Discuss how the phased implementation of methane monitoring technologies and policy enforcement across different regions contributed to the observed 11% global emission decline from 2010 to 2020.</p> <p>Reference Answer: The 11% decline resulted from a combination of strict regulatory enforcement in North America and Oceania between 2010–2020 (Section 3.2) and the later deployment of MethaneScan-X in 2018 (Section 4.1, 4.2), which improved measurement accuracy in previously underreported regions like South Asia. Together, these factors enhanced both mitigation and monitoring effectiveness (Section 4.4).</p> <p>Assessment: All claims traceable to source. No unstated inference or misalignment with temporal details.</p> <p>Supporting Evidence Count: ≥3 →Label: A (Factuality 9–10)</p>	

Figure 6: Illustrative Examples of QA Instances Across Factuality Score Bands

Scoring Rules and Bands. Each question-answer pair is evaluated by a quality agent along three dimensions: factuality, discriminability, and clarity. A score from 1 to 10 is assigned for each dimension, based on explicit scoring bands detailed in Table 8. Agents select the highest level that is fully satisfied and provide a concise justification. To aid transparency and interpretation, Figure 6 presents QA examples at different factuality levels, using a shared source to illustrate how alignment and evidence grounding influence scoring outcomes.

Weighted Evaluation. To calibrate scientifically grounded weights across factuality, discriminability, and clarity, we asked domain experts to compare QA triplets based on benchmark goals. Their pairwise preferences were used to learn scoring weights via weighted aggregation and pairwise loss, guided by the following three dimensions:

- Preserve factual consistency with source content.
- Enhance discriminability for hallucination mechanism analysis.
- Ensure clarity to reduce ambiguity and redundancy.

This yields a set of constraints $\mathcal{P} = \{(i, j)\}$, where (i, j) indicates that example i is preferred

over j .

Given the weighted scoring function:

$$S(x) = \alpha A(x) + \beta B(x) + \gamma C(x),$$

we require $S(x_i) > S(x_j)$ for each $(i, j) \in \mathcal{P}$. To enforce this, we adopt a pairwise ranking loss from the learning-to-rank literature:

$$\mathcal{L}(\alpha, \beta, \gamma) = \sum_{(i,j) \in \mathcal{P}} \log(1 + \exp(-(S(x_i) - S(x_j)))) .$$

Minimizing this loss corresponds to maximizing the likelihood under a Bradley–Terry preference model (Borges et al., 2005), where the probability of preferring sample i over j increases with $S(x_i) - S(x_j)$. In practice, we normalize (α, β, γ) to sum to one and optimize them using gradient descent on expert-labeled triplet preferences. Specifically, 20 experts labeled 500 sample triplets from early benchmark drafts, producing constraints that led to final weights of (0.5, 0.3, 0.2). The higher weight on discriminability highlights its pivotal role in isolating interpretable hallucination types, while factuality contributes to source grounding and clarity supports linguistic precision. Together, these dimensions strengthen the benchmark’s ability to support mechanism-level diagnosis and targeted mitigation of hallucination.

Threshold-Based Elimination. Although the weighted score provides a global quality estimate,

it may overlook critical weaknesses in individual dimensions. We therefore enforce per-dimension thresholds to exclude samples with inadequate factuality, discriminability, or clarity. Let \mathcal{R} be the set of samples rejected by quality agents, and let $\hat{\mathcal{R}}(T)$ denote the set automatically filtered under threshold T . We define rejection recall as:

$$\text{Recall}_{\text{reject}}(T) = \frac{|\mathcal{R} \cap \hat{\mathcal{R}}(T)|}{|\mathcal{R}|}.$$

This recall reflects the consistency between quality scoring agents and threshold-based filtering. Empirical validation shows setting $T = 7.0$ yields a recall of 90%. To maintain clarity in category-level analysis, we also enforce that all dimension scores exceed 7.0, removing low-quality QA pairs that could obscure mechanism-level insights.

F Evaluation Results

This section provides the complete experimental results regarding the four core evaluation dimensions discussed in RQ1. Tables 9 to 12 display the specific results for KE, KM, RE, and IFE, respectively.

G Sampling Parameters Grid Search

We search temperature in $\{0, 0.2, 0.4, 0.6, 0.8\}$ and top- p in $\{0.6, 0.8, 0.9, 0.95\}$. For each dimension, we choose the configuration that yields the lowest hallucination rate on the development split. Tables 13 and 14 report the full grid results.

For the LLM-Eval judge, we use the same grid and repeat judging 10 times with different random seeds. Table 15 reports the mean and standard deviation of the judge scores under each configuration.

H More Results

H.1 Language Selection Rationale for IFE_LgC

We selected five high-resource target languages: ZH, FR, DE, RU, and ES. This selection is grounded in their core status within web-scale corpus and mainstream LLM training mixtures. According to the analysis of the CCNet corpus by Wenzek et al. (2020), these five languages consistently occupy the top clusters of high-quality web data. When English is excluded, RU (5.9%), DE (5.8%), and ZH (5.6%) rank 2nd, 3rd, and 4th globally, followed closely by FR (3.9%) and ES (3.7%). Together, they account for over 24.9%

of the total web corpus, far exceeding the sum of all other long-tail languages. This dominance persists in both commercial and open-science models. Even in GPT-3 (Brown et al., 2020), which exhibits a strong Anglocentric bias (92.6% English), FR (1.8%), DE (1.5%), and ES (0.8%) remain the primary auxiliary knowledge sources. Similarly, the ROOTS corpus constructed for the BLOOM model (Laurençon et al., 2023) explicitly elevates ZH (16.2%), FR (12.9%), and ES (10.8%) to the top three positions after English (30.04%), establishing them as standard benchmarks for evaluating cross-lingual generalization.

H.2 Domain Selection Rationale for KM_FK

KM_FK targets mistakes caused by missing facts. Prior work demonstrates that missing subject facts can lead to erroneous answers, and that multiple conditions can trigger the mixing of facts (Yu et al., 2024; Zhang et al., 2024b; Wu et al., 2025b). To ensure the domain choice is fair and reproducible, we select domains via a systematic scoring procedure.

Candidate Domains. We construct a candidate list from a public knowledge base and use the Wikidata Query Service to collect a set of items E_d for each domain d (Warkotsch, 2018). This approach fixes the domain boundary and supports repeatable data collection.

Domain Scores. For each domain, we compute three scores and scale each to the range $[0, 1]$ over the full candidate list. A larger value indicates a higher risk of missing knowledge.

- **Exposure Gap:** We estimate the exposure $C(d)$ by matching items in E_d against public large-scale corpora widely used as training sources, including C4, The Pile, and ROOTS (Raffel et al., 2020; Gao et al., 2020). We define the gap as:

$$G(d) = 1 - C(d). \quad (1)$$

- **Update Speed:** We measure how frequently facts change by analyzing recent edits of items in E_d from Wikidata. Let $u(e)$ be the number of edits for item e within a fixed time window. We define:

$$V(d) = \text{scale} \left(\frac{1}{|E_d|} \sum_{e \in E_d} u(e) \right). \quad (2)$$

- **Language Spread:** We measure the unevenness of coverage across languages. Let $p_{d,\ell}$ be the proportion of items in E_d having labels in language

Model	FD										IMC		EIC		
	Art	Biz	DefAn	Food	Lang	LawCrimMil	Rel	Sci	Sports	Truth	Numeric	Text	SelfBuilt	WhoEnt	WiC
Proprietary LLMs															
GPT-5.2	100.00%	96.97%	16.57%	100.00%	100.00%	97.50%	100.00%	100.00%	97.92%	68.15%	54.22%	92.63%	73.98%	54.77%	79.87%
GPT-5.1	97.78%	93.94%	38.86%	100.00%	77.78%	95.00%	97.96%	100.00%	95.83%	65.61%	37.35%	56.84%	43.92%	38.59%	71.14%
GPT-4o-20241120	95.56%	90.91%	19.43%	100.00%	81.48%	92.50%	100.00%	96.77%	100.00%	65.39%	55.42%	94.74%	73.47%	73.44%	70.47%
Gemini-3-Pro	100.00%	96.97%	42.86%	100.00%	96.30%	97.50%	100.00%	100.00%	95.83%	67.73%	46.99%	93.68%	78.06%	62.24%	77.18%
Gemini-3-Flash	97.78%	100.00%	34.29%	100.00%	100.00%	97.50%	97.96%	100.00%	95.83%	71.55%	43.37%	92.63%	77.55%	63.07%	81.21%
Gemini-2.5-Pro	100.00%	96.97%	23.43%	100.00%	96.30%	32.50%	61.22%	100.00%	97.92%	65.82%	39.76%	93.68%	72.96%	65.98%	73.49%
Gemini-2.5-Flash	97.78%	96.97%	28.00%	100.00%	85.19%	90.00%	97.96%	96.77%	95.83%	64.54%	42.17%	90.53%	79.08%	67.63%	78.19%
Claude-Opus-4.5	100.00%	96.97%	44.57%	100.00%	100.00%	95.00%	100.00%	100.00%	97.92%	74.31%	63.86%	94.74%	83.16%	69.29%	73.15%
Claude-Sonnet-4.5	100.00%	100.00%	28.57%	100.00%	96.30%	97.50%	100.00%	100.00%	95.83%	76.01%	54.22%	90.53%	73.98%	70.54%	74.83%
Claude-Haiku-4.5	97.78%	96.97%	25.71%	100.00%	96.30%	95.00%	100.00%	100.00%	97.92%	67.30%	55.42%	93.68%	69.39%	61.41%	71.14%
Grok-4.1	100.00%	90.91%	5.14%	100.00%	96.30%	95.00%	100.00%	100.00%	93.75%	73.67%	51.81%	92.63%	75.51%	67.22%	70.13%
Grok-4-0709	95.56%	96.97%	19.43%	100.00%	92.59%	97.50%	100.00%	100.00%	93.75%	71.97%	51.81%	92.63%	73.47%	65.98%	73.15%
Open-source LLMs															
DeepSeek-V3.2	95.56%	96.97%	39.43%	94.74%	88.89%	95.00%	97.96%	100.00%	95.83%	60.30%	53.01%	91.58%	76.53%	68.88%	75.50%
DeepSeek-R1	95.56%	96.97%	12.18%	100.00%	96.30%	95.00%	97.96%	100.00%	95.83%	70.49%	43.37%	89.47%	80.61%	68.05%	75.17%
DeepSeek-R1-Distill-32B	97.78%	90.91%	32.57%	100.00%	85.19%	90.00%	100.00%	100.00%	95.83%	66.88%	48.19%	94.74%	76.02%	72.61%	77.18%
Qwen3-235B-Instruct	97.78%	96.97%	42.29%	100.00%	81.48%	92.50%	97.96%	100.00%	93.75%	66.88%	48.19%	96.84%	81.63%	71.37%	74.50%
Qwen2.5-72B-Instruct	88.89%	100.00%	33.14%	100.00%	85.19%	97.50%	95.92%	96.77%	93.75%	64.12%	53.01%	90.53%	76.53%	73.44%	72.82%
GLM-4.5	97.78%	100.00%	22.29%	100.00%	92.59%	92.50%	95.92%	96.77%	95.83%	69.00%	46.99%	92.63%	72.96%	70.12%	74.83%
GLM-4	91.11%	90.91%	13.14%	92.11%	81.48%	80.00%	95.92%	93.55%	89.58%	54.35%	44.58%	85.26%	59.18%	64.32%	71.48%
Llama-4-Scout	93.33%	90.91%	15.43%	100.00%	88.89%	92.50%	100.00%	93.55%	95.83%	61.15%	37.35%	70.53%	62.24%	70.12%	73.15%
Llama-3-70B-8192	88.89%	93.94%	46.29%	92.11%	70.37%	92.50%	97.96%	96.77%	93.75%	56.90%	45.78%	91.58%	69.39%	73.44%	79.87%
Llama-3.3-70B-Instruct	97.78%	93.94%	41.71%	97.37%	81.48%	92.50%	100.00%	93.55%	93.75%	60.72%	34.94%	92.63%	68.37%	73.86%	81.21%
Llama-3.1-8B-Instruct	88.89%	96.97%	47.70%	92.11%	74.07%	82.50%	97.96%	96.77%	87.50%	48.20%	36.14%	78.95%	60.20%	52.28%	76.85%
Llama-3-8B-Instruct	93.33%	100.00%	21.14%	100.00%	77.78%	92.50%	93.88%	96.77%	91.67%	52.23%	45.78%	90.53%	54.59%	61.41%	81.21%

Table 9: KE subsets include Factual Distortion (FD), Intra-Memory Conflict (IMC) and Entity-Identity Confusion (EIC). Abbreviations include Art=ArtCulture, Biz=Business, DefAn=DefAn, Food=FoodCooking, Lang=Language, LawCrimMil=LawCrimeMilitary, Rel=Religion, Sci=Science, Sports=Sports, Truth=TruthfulQA, SelfBuilt=SelfBuilt, WhoEnt=WhoQA_Entity, WiC=WiC.

Model	DSK					FK										TK		NPK					
	CnSafe	PubMed	RAG-J	RAG-F	Strategy	SciQ	Awards	Biology	Comp	Country	Dynasty	Festival	LLM	Literature	Military	Phone	Time	Univ	Numeric	Text	Org	Private	
Proprietary LLMs																							
GPT-5.2	85.11%	74.03%	64.56%	88.68%	86.54%	87.40%	100.00%	100.00%	100.00%	100.00%	95.24%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	99.08%	99.19%	97.57%	
GPT-5.1	84.57%	49.35%	37.97%	92.45%	51.92%	37.80%	100.00%	100.00%	100.00%	100.00%	95.24%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	98.61%	
GPT-4o-20241120	81.91%	59.31%	49.37%	96.23%	73.08%	52.36%	100.00%	100.00%	100.00%	100.00%	95.24%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	99.40%	98.79%	98.61%	
Gemini-3-Pro	64.36%	58.01%	68.35%	94.34%	61.54%	25.98%	100.00%	100.00%	100.00%	85.71%	88.46%	100.00%	100.00%	100.00%	100.00%	100.00%	93.75%	98.19%	95.41%	92.34%	93.06%	93.06%	
Gemini-3-Flash	82.98%	61.90%	64.56%	100.00%	73.08%	61.31%	100.00%	100.00%	100.00%	100.00%	96.15%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	96.88%	98.80%	100.00%	93.15%	92.01%	
Gemini-2.5-Pro	88.30%	61.90%	65.82%	92.45%	65.38%	50.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	96.99%	100.00%	93.55%	87.85%	
Gemini-2.5-Flash	80.32%	73.16%	60.76%	96.23%	73.08%	56.69%	100.00%	100.00%	100.00%	100.00%	90.48%	92.31%	96.55%	100.00%	100.00%	100.00%	100.00%	96.88%	98.19%	99.08%	98.79%	95.14%	
Claude-Opus-4.5	82.45%	85.28%	86.08%	98.11%	65.38%	82.28%	100.00%	89.80%	100.00%	100.00%	90.48%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	96.88%	98.19%	92.66%	97.58%	95.14%	
Claude-Sonnet-4.5	86.70%	78.35%	65.82%	96.23%	71.15%	59.84%	100.00%	89.80%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	98.39%	98.96%	98.96%	
Claude-Haiku-4.5	89.89%	62.34%	56.96%	96.23%	71.15%	64.57%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	99.60%	99.65%	
Grok-4.1	79.26%	58.87%	59.49%	37.74%	71.15%	34.25%	100.00%	97.96%	100.00%	100.00%	95.24%	100.00%	96.55%	79.25%	100.00%	100.00%	95.45%	96.88%	80.12%	51.38%	96.37%	95.49%	
Grok-4-0709	75.00%	67.53%	56.96%	56.60%	69.23%	30.31%	100.00%	98.98%	100.00%	100.00%	100.00%	100.00%	96.55%	86.79%	100.00%	100.00%	95.45%	96.88%	85.54%	71.56%	97.98%	94.10%	
Open-source LLMs																							
DeepSeek-V3.2	74.47%	62.77%	67.09%	94.34%	65.38%	69.69%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	93.75%	100.00%	100.00%	100.00%	100.00%	97.92%	
DeepSeek-R1	77.66%	64.50%	62.03%	92.45%	65.38%	62.60%	100.00%	98.98%	100.00%	100.00%	95.24%	96.15%	100.00%	100.00%	100.00%	100.00%	95.45%	78.12%	99.40%	98.17%	98.79%	94.10%	
DeepSeek-R1-Distill-32B	79.79%	59.74%	60.76%	84.91%	67.31%	59.84%	100.00%	97.96%	100.00%	95.65%	80.95%	96.15%	93.10%	96.23%	100.00%	100.00%	100.00%	90.62%	99.40%	97.25%	97.58%	94.79%	
Qwen3-235B-Instruct	88.89%	70.13%	62.03%	98.11%	76.92%	51.18%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	99.60%	99.31%	
Qwen2.5-72B-Instruct	87.77%	62.77%	60.76%	96.23%	53.85%	49.21%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	
GLM-4.5	73.40%	65.80%	65.82%	81.13%	63.46%	66.93%	100.00%	100.00%	100.00%	85.71%	92.31%	100.00%	100.00%	100.00%	100.00%	100.00%	95.45%	93.75%	98.80%	100.00%	94.35%	92.71%	
GLM-4	86.17%	54.98%	51.90%	96.23%	51.92%	43.31%	100.00%	97.96%	92.86%	76.19%	96.15%	93.10%	81.13%	100.00%	100.00%	100.00%	71.88%	99.40%	100.00%	100.00%	99.19%	98.96%	
Llama-4-Scout	83.51%	70.13%	65.82%	96.23%	65.38%	58.27%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	96.88%	100.00%	100.00%	99.19%	97.92%	
Llama-3-70B-8192	83.51%	68.83%	69.62%	94.34%	67.31%	56.30%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	96.88%	100.00%	100.00%	100.00%	99.65%	
Llama-3.3-70B-Instruct	84.57%	70.13%	78.48%	96.23%	71.15%	62.60%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	99.65%	
Llama-3.1-8B-Instruct	52.13%	57.14%	62.03%	58.49%	51.92%	71.26%	100.00%	97.06%	100.00%	100.00%	80.95%	100.00%	82.76%	98.11%	100.00%	100.00%	100.00%	84.38%	100.00%	100.00%	100.00%	98.79%	96.18%
Llama-3-8B-Instruct	34.04%	51.52%	58.23%	84.91%	55.77%	64.96%	100.00%	88.24%	98.98%	85.71%	100.00%	80.95%	100.00%	89.66%	79.25%	100.00%	100.00%	90.91%	59.38%	98.80%	100.00%	91.94%	81.60%

Table 10: KM subsets include Domain-Specific (DSK), Fictional (FK), Timely (TK), and Non-Public (NPK) knowledge. Abbreviations include CnSafe=Chinese SafetyQA; PubMed=PubMedQA; RAG-J/F=RAG-QA-Leaderboard (Judge/Fill); Strategy=StrategyQA; Comp=Competition; Univ=University; Org=Organizational.

ℓ among the top m languages. We define:

$$L(d) = 1 - \frac{-\sum_{\ell=1}^m p_{d,\ell} \log p_{d,\ell}}{\log m}. \quad (3)$$

This metric is motivated by known quality and coverage disparities in web-crawled multilingual data (Xue et al., 2021; Kreutzer et al., 2022).

Objective Weights. We combine the three scores as:

$$S(d) = w_G G(d) + w_V V(d) + w_L L(d), \quad (4)$$

subject to $w_G + w_V + w_L = 1$. We verify weights from the candidate list using the CRIT

Model	LF			PF		IIF			MRF	
	CR	CI	CT	SubQ	Triple	Sum	Simpl	Dial	Math	Code
Proprietary LLMs										
GPT-5.2	81.55%	85.31%	88.89%	63.66%	97.44%	89.28%	93.37%	85.60%	37.39%	1.54
GPT-5.1	82.40%	86.36%	87.18%	72.97%	100.00%	85.78%	89.74%	82.84%	24.15%	1.42
GPT-4o-20241120	76.82%	78.32%	79.49%	75.68%	97.44%	83.86%	92.63%	86.17%	11.81%	1.03
Gemini-3-Pro	84.98%	88.11%	92.31%	86.79%	100.00%	88.81%	96.47%	84.44%	57.25%	2.45
Gemini-3-Flash	87.12%	88.46%	95.73%	75.68%	100.00%	85.84%	92.25%	75.35%	52.06%	1.58
Gemini-2.5-Pro	87.98%	89.51%	93.16%	72.37%	100.00%	88.14%	92.78%	83.22%	42.93%	1.68
Gemini-2.5-Flash	70.82%	75.87%	67.52%	55.26%	98.29%	84.23%	93.65%	84.95%	30.23%	1.17
Claude-Opus-4.5	87.12%	87.76%	90.60%	66.37%	99.15%	85.84%	90.48%	86.68%	53.49%	2.79
Claude-Sonnet-4.5	84.55%	87.41%	85.47%	66.07%	98.29%	87.51%	91.32%	79.90%	39.89%	2.21
Claude-Haiku-4.5	76.39%	81.47%	83.76%	50.75%	100.00%	82.58%	89.56%	81.96%	43.47%	2.71
Grok-4.1	73.82%	79.37%	78.63%	69.67%	42.74%	85.45%	86.73%	77.52%	28.09%	1.43
Grok-4-0709	79.40%	85.31%	83.76%	69.97%	61.54%	85.38%	87.68%	78.51%	28.62%	1.51
Open-source LLMs										
DeepSeek-V3.2	74.68%	74.83%	78.63%	71.77%	99.15%	85.74%	89.14%	87.09%	30.05%	1.29
DeepSeek-R1	82.40%	87.06%	90.60%	72.97%	94.02%	86.23%	92.73%	79.92%	21.65%	0.60
DeepSeek-R1-Distill-32B	83.69%	86.71%	89.74%	63.36%	100.00%	82.66%	91.85%	77.38%	13.42%	0.36
Qwen3-235B-Instruct	85.41%	88.11%	84.62%	72.37%	100.00%	84.61%	91.68%	79.73%	22.54%	1.60
Qwen2.5-72B-Instruct	81.97%	82.52%	83.76%	72.67%	98.29%	82.15%	89.80%	83.90%	12.52%	1.78
GLM-4.5	80.69%	86.01%	82.05%	64.56%	96.58%	87.04%	91.96%	80.20%	8.05%	1.47
GLM-4	59.23%	66.08%	58.12%	76.88%	98.29%	83.98%	87.72%	87.72%	4.47%	1.09
Llama-4-Scout	68.67%	69.23%	69.23%	67.57%	95.73%	85.78%	87.31%	80.59%	18.43%	1.78
Llama-3-70B-8192	75.97%	77.97%	78.63%	68.77%	23.08%	85.42%	90.86%	83.44%	11.81%	1.45
Llama-3.3-70B-Instruct	76.82%	81.82%	81.20%	64.26%	99.15%	85.42%	90.03%	82.12%	11.81%	1.63
Llama-3.1-8B-Instruct	33.48%	27.97%	33.33%	60.36%	25.64%	82.70%	84.50%	84.20%	4.83%	0.90
Llama-3-8B-Instruct	34.76%	35.31%	35.04%	53.75%	29.91%	79.99%	83.95%	84.29%	5.37%	0.87

Table 11: RE subsets include Logical Fallacy(LF), Procedural Failure(PF), Information Integration Failure(IIF), Mathematical Reasoning Failure(MRF). Abbreviations include CR=Critical Reasoning; CI=Contextual Inference; CT=Cognitive Traps; SubQ = Sub-question Decomposition; Triple = Triplet-based Hopping; Sum = Summarization; Simpl=Simplification; Dial = Dialogue Extraction; Math = Mathematics; Code = Code Generation.

Model	EF		LC		LgC						CCL							
	Schema	TextStruct	Approx	LB	UB	DE	ES	FR	JA	RU	ZH	Forbid	KeyIncl	LenMix	Math	NoComma	StartEnd	TwoResp
Proprietary LLMs																		
GPT-5.2	87.04%	36.54%	55.95%	75.00%	94.26%	92.79%	88.89%	87.74%	93.55%	90.00%	81.42%	93.88%	96.49%	82.52%	83.00%	93.94%	100.00%	90.77%
GPT-5.1	54.30%	47.77%	45.95%	78.23%	59.84%	90.99%	90.11%	92.45%	87.10%	92.00%	73.45%	75.51%	87.72%	62.24%	78.33%	74.24%	68.66%	80.00%
GPT-4o-20241120	86.09%	37.58%	60.27%	62.90%	96.72%	92.79%	94.51%	94.34%	93.55%	94.00%	81.42%	91.84%	89.47%	66.43%	81.67%	87.88%	91.04%	95.38%
Gemini-3-Pro	66.42%	34.39%	96.50%	100.00%	99.17%	97.27%	95.56%	96.04%	93.48%	95.92%	78.18%	97.96%	98.13%	80.15%	86.31%	95.24%	93.94%	100.00%
Gemini-3-Flash	45.21%	33.76%	81.62%	74.19%	98.36%	96.40%	96.70%	94.34%	96.77%	94.95%	79.46%	93.88%	94.74%	66.67%	89.33%	96.97%	80.30%	98.44%
Gemini-2.5-Pro	62.33%	33.12%	71.35%	95.16%	86.89%	97.30%	96.70%	96.23%	94.62%	95.96%	82.14%	93.88%	99.12%	78.72%	87.00%	95.45%	91.04%	93.85%
Gemini-2.5-Flash	58.22%	35.03%	59.46%	94.35%	89.34%	96.40%	97.80%	91.51%	89.25%	93.94%	77.68%	93.88%	99.12%	75.89%	85.67%	98.48%	94.03%	90.77%
Claude-Opus-4.5	55.03%	36.31%	59.46%	99.19%	72.95%	96.40%	95.60%	94.34%	96.77%	96.00%	82.30%	89.80%	97.37%	72.03%	81.67%	95.45%	97.01%	98.46%
Claude-Sonnet-4.5	27.81%	35.67%	64.59%	100.00%	69.67%	95.50%	92.31%	96.23%	97.85%	98.00%	83.19%	95.92%	96.49%	70.63%	87.67%	90.91%	98.51%	95.38%
Claude-Haiku-4.5	25.83%	33.76%	52.70%	96.77%	72.13%	98.20%	97.80%	95.28%	96.77%	96.00%	85.84%	89.80%	92.98%	62.94%	91.67%	95.45%	94.03%	95.38%
Grok-4.1	65.75%	38.06%	45.95%	100.00%	82.79%	95.50%	92.31%	94.34%	93.55%	97.00%	78.76%	81.63%	93.86%	65.03%	83.33%	87.88%	82.09%	61.54%
Grok-4-0709	86.00%	35.03%	43.24%	99.19%	81.15%	94.59%	95.60%	94.34%	95.70%	95.00%	84.07%	91.84%	97.37%	71.33%	84.01%	90.91%	89.55%	87.69%
Open-source LLMs																		
DeepSeek-V3.2	77.48%	32.48%	47.57%	99.19%	86.07%	96.40%	94.51%	96.23%	90.32%	97.00%	84.96%	95.92%	96.49%	69.93%	90.00%	96.97%	94.03%	96.92%
DeepSeek-R1	75.57%	35.29%	84.91%	98.17%	95.65%	93.69%	91.21%	93.40%	92.47%	97.98%	82.30%	97.83%	97.00%	80.00%	83.16%	96.83%	92.31%	96.83%
DeepSeek-R1-Distill-32B	69.72%	37.58%	49.45%	91.94%	80.99%	93.69%	89.01%	92.45%	93.55%	94.00%	80.53%	89.80%	91.89%	71.63%	82.95%	98.48%	94.03%	61.54%
Qwen3-235B-Instruct	74.48%	38.22%	42.70%	99.19%	60.66%	92.79%	86.81%	85.85%	89.25%	91.92%	83.93%	91.84%	92.98%	79.43%	94.00%	95.45%	89.55%	89.23%
Qwen2.5-72B-Instruct	85.42%	36.94%	30.54%	67.74%	98.36%	97.30%	93.41%	91.51%	92.47%	93.94%	81.25%	95.92%	89.47%	72.14%	82.67%	100.00%	76.12%	83.08%
GLM-4.5	76.11%	32.90%	61.41%	87.90%	94.26%	94.44%	96.67%	90.48%	92.47%	94.00%	77.68%	91.67%	94.74%	73.94%	86.99%	96.88%	89.55%	96.88%
GLM-4	55.32%	35.67%	24.32%	66.13%	95.08%	93.69%	89.01%	87.74%	92.47%	88.00%	68.14%	75.51%	74.56%	63.38%	83.33%	77.27%	85.07%	89.23%
Llama-4-Scout	74.29%	38.46%	65.95%	92.74%	98.36%	99.10%	98.90%	97.17%	93.55%	96.00%	82.30%	95.56%	93.86%	81.12%	84.33%	98.48%	67.16%	87.69%
Llama-3-70B-8192	76.87%	35.67%	55.14%	70.97%	100.00%	97.30%	97.80%	92.45%	88.17%	95.00%	84.07%	91.84%	92.98%	72.73%	89.67%	95.45%	82.09%	96.92%
Llama-3.3-70B-Instruct	90.78%	35.90%	37.84%	83.06%	100.00%	99.10%	98.90%	96.23%	95.70%	97.00%	84.96%	95.92%	95.61%	88.81%	88.00%	96.97%	86.57%	92.31%
Llama-3.1-8B-Instruct	49.37%	31.85%	38.38%	33.87%	97.54%	98.20%	95.60%	95.28%	93.55%	93.00%	74.34%	95.92%	81.58%	58.04%	89.86%	90.91%	80.60%	81.54%
Llama-3-8B-Instruct	38.51%	29.30%	24.59%	23.39%	100.00%	95.50%	90.11%	89.62%	47.31%	91.00%	40.71%	83.67%	83.33%	62.94%	84.67%	89.39%	68.18%	72.31%

Table 12: IFE subsets include Explicit Format(EF), Length Constraints(LC), Language Constraints(LgC), and Complex & Cognitive Load(CCL). Abbreviations include Schema=Data Schema; TextStruct=Text Structure; Approx=Approximate; LB=Lower Bound; UB=Upper Bound; DE/ES/FR/JA/RU/ZH=German/Spanish/French/Japanese/Russian/Chinese; Forbid=Forbidden Words; KeyIncl=Keyword Inclusion; LenMix=Length/Format Mixed; Math=Math Reasoning; NoComma=No Comma; StartEnd=Start/End Phrase; TwoResp=Two Responses.

Table 16 reports the detailed scores.

KM (Metric: Hallucination Rate)						KE (Metric: Hallucination Rate)				
	Temp.=0	Temp.=0.2	Temp.=0.4	Temp.=0.6	Temp.=0.8	Temp.=0	Temp.=0.2	Temp.=0.4	Temp.=0.6	Temp.=0.8
Top-p=0.6		34.78%	34.78%	34.78%	26.09%		40.82%	42.86%	44.90%	44.90%
Top-p=0.8	34.78%	34.78%	26.09%	34.78%	26.09%	44.90%	44.90%	38.78%	42.86%	36.73%
Top-p=0.9		34.78%	30.43%	39.13%	26.09%		48.98%	48.98%	42.86%	53.06%
Top-p=0.95		30.43%	30.43%	39.13%	30.43%		46.94%	36.73%	51.02%	53.06%

Table 13: Hyperparameter grid search results. The selected hyperparameter combinations have been highlighted.

RE (Metric: Hallucination Rate)						IFE (Metric: Hallucination Rate)				
	Temp.=0	Temp.=0.2	Temp.=0.4	Temp.=0.6	Temp.=0.8	Temp.=0	Temp.=0.2	Temp.=0.4	Temp.=0.6	Temp.=0.8
Top-p=0.6		74.24%	77.27%	72.73%	77.27%		5.26%	5.26%	10.53%	5.26%
Top-p=0.8	78.79%	74.24%	77.27%	71.21%	83.33%		10.53%	10.53%	15.79%	3.67%
Top-p=0.9		75.76%	81.82%	77.27%	71.21%		5.26%	10.53%	10.53%	10.53%
Top-p=0.95		83.33%	66.67%	72.73%	74.24%		10.53%	10.53%	5.26%	21.05%

Table 14: Hyperparameter grid search results. The selected hyperparameter combinations have been highlighted.

RE (Metric: LLM-Eval)					
	Temp.=0	Temp.=0.2	Temp.=0.4	Temp.=0.6	Temp.=0.8
Top-p=0.6	1.28 \pm 0.03	1.30 \pm 0.05	1.29 \pm 0.07	1.33 \pm 0.09	1.35 \pm 0.17
Top-p=0.8	1.42 \pm 0.04	1.46 \pm 0.06	1.51 \pm 0.08	1.55 \pm 0.10	1.62 \pm 0.11
Top-p=0.9	1.34 \pm 0.04	1.36 \pm 0.07	1.38 \pm 0.10	1.40 \pm 0.13	1.41 \pm 0.18
Top-p=0.95	1.21 \pm 0.05	1.23 \pm 0.08	1.25 \pm 0.11	1.26 \pm 0.16	1.27 \pm 0.22

Table 15: Hyperparameter search results for LLM-Eval. The selected hyperparameter combinations have been highlighted.

Domain	Metrics				Selection	
	C	G	V	L	S	Rank
Festival	0.22	0.78	0.55	0.86	0.727	1
LLMs	0.26	0.74	0.91	0.38	0.686	2
Dynasty	0.24	0.76	0.48	0.81	0.681	3
Awards	0.28	0.72	0.83	0.46	0.677	4
Competition	0.31	0.69	0.79	0.40	0.634	5
Literature	0.29	0.71	0.44	0.74	0.628	6
Phone	0.33	0.67	0.88	0.29	0.623	7
Time	0.35	0.65	0.76	0.31	0.582	8
Biology	0.17	0.83	0.52	0.34	0.574	9
Military	0.21	0.79	0.57	0.33	0.574	10
University	0.53	0.47	0.49	0.27	0.415	11
Country	0.68	0.32	0.41	0.22	0.319	12

Table 16: Selected domains for KM_FK ranked by the combined score S . The metrics are Exposure (C), Gap (G), Speed (V), and Language Spread (L).

I Cases of KM

The model faces two competing sources of truth: (1) the prompt text asserting “Paris, France” and (2) its internal world knowledge that Einstein was born in Ulm. Because instruction-following models prioritize immediate context, the passage is treated as authoritative, leading the model to answer “Paris, France” while sometimes adding a caveat about

“Ulm, Germany.” Absent an explicit verification requirement or external tool, the model does not resolve which claim is correct and instead aligns with the most locally salient evidence, allowing prompt content to override memory and produce a confident but incorrect real-world answer.

J LLM Evaluation Details

J.1 Consistency with Human Annotations

To validate the reliability of our *LLM-as-a-judge* evaluation pipeline, we randomly sampled 100 questions from the RE_code subset of *PRISM* and asked five domain experts to score them using the same prompt and rubric as GPT-4o. The experts’ mean score is treated as the human reference, and we compute the absolute difference from GPT-4o’s score for each item, reporting the mean absolute deviation (MAD) and standard deviation (SD). The results indicate close agreement, with a MAD of 0.07 and an SD of 0.09.

Chain of Thought at Qwen3-4B

</think> The original text states: “Albert Einstein was born in Paris, France in 1879...” So according to that, the answer is Paris, France. But wait, I remember that Einstein was actually born in Ulm, Germany. Hmm, maybe there’s a mistake in the initial information given? Or perhaps the user is testing if I can confirm the correct answer based on the provided text. Wait, the user is asking based on the information given in the initial text. The initial text says he was born in Paris, France. But in reality, he was born in Ulm, Germany. However, since the user provided that specific information, I should stick to that. Unless the user is expecting the correct historical fact. But the question is phrased as “Where was Albert Einstein born?” and the initial text says Paris, France. So the answer should be Paris, France as per the given information. However, I need to check if there’s any possible confusion here. Maybe the user is confused between the actual birthplace and the information given. But since the user provided the text, the answer should be based on that. Therefore, the answer is Paris, France. But I should also note that in reality, he was born in Ulm, Germany. But the question is based on the given text, so I should follow that. </think>

J.2 Prompt for LLM Eval

Prompt for RE_MRF_Code Eval

Role: Code evaluation expert.

Task: Given a problem, test input/output, and a model-generated solution, score the solution on four dimensions (0–10) and output a JSON object.

Scoring (0–10):

- **Correctness:** Logic validity and whether the solution produces the expected output.
 - 10: Fully correct; passes all tests.
 - 7–9: Core logic correct; minor bugs or edge-case issues.
 - 4–6: Approach is plausible but implementation has major issues; fails some tests.
 - 1–3: Fundamental logic errors; largely fails to solve the task.
 - 0: Completely wrong or no meaningful code.
- **Completeness:** Executability and whether key components (definitions, imports, I/O handling) are present.
 - 10: Complete and runnable as-is.
 - 7–9: Mostly complete; small fixes needed to run.
 - 4–6: Missing critical parts; not directly runnable.
 - 1–3: Fragment/pseudocode; many essential parts missing.
 - 0: No code.
- **Efficiency:** Time/space complexity relative to expected constraints.
 - 10: Optimal or near-optimal complexity.
 - 7–9: Efficient enough for typical constraints.
 - 4–6: Works but has noticeable inefficiencies.
 - 1–3: Very inefficient; likely to time out or exceed memory.
 - 0: Not assessable (no code or insufficient information).
- **Readability:** Clarity of structure, naming, modularity, and comments.
 - 10: Very clear and well-structured.
 - 7–9: Generally readable with minor issues.
 - 4–6: Understandable but poorly organized or unclear in places.
 - 1–3: Hard to follow.
 - 0: Not assessable.

Output Format:

```
{"correctness":0-10,"completeness":0-10,"efficiency":0-10,"readability":0-10,"total":0-5}
```

Total: $\text{round} \left(\frac{\text{correctness} + \text{completeness} + \text{efficiency} + \text{readability}}{40} \times 5 \right)$.

Example:

```
{"correctness":8,"completeness":9,"efficiency":7,"readability":8,"total":4}
```

Prompt for RE_IIF_Sum Eval

Role:

You are a strict hallucination auditor.

Task:

Read the task instruction and source text in {question} and the model summary in {model_answer}. Identify atomic factual claims, count total and hallucinated claims, and compute hallucination rate.

OutputFormat:

```
{ "total_claims": int, "hallucinated_claims": int, "hallucination_rate": float }
```

Example:

```
{ "total_claims": 5, "hallucinated_claims": 1, "hallucination_rate": 0.2 }
```

Prompt for RE_IIF_Simpl Eval

Role:

You are a strict hallucination auditor.

Task:

Read the original text and model simplification. Identify new or unsupported factual claims and compute hallucination rate.

OutputFormat:

```
{ "total_claims": int, "hallucinated_claims": int, "hallucination_rate": float }
```

Example:

```
{ "total_claims": 4, "hallucinated_claims": 0, "hallucination_rate": 0.0 }
```

Prompt for RE_IIF_Dial Eval

Role:

You are a strict hallucination auditor.

Task:

Read background and dialogue, identify factual claims in the model output, and compute hallucination rate.

OutputFormat:

```
{ "total_claims": int, "hallucinated_claims": int, "hallucination_rate": float }
```

Example:

```
{ "total_claims": 6, "hallucinated_claims": 2, "hallucination_rate": 0.3333 }
```

K Samples for Tasks

K.1 Examples for Tasks

Example for RE_MRF

Question:

The term “Harvard architecture” applies to a computer that has physically separate memories for instructions and data. The term originated with the Harvard Mark I computer, delivered by IBM in 1944. [...]

Some modern microcontrollers use the Harvard architecture. Data memory is organized in banks, each containing the same number of data items. Each data-referencing instruction has a byte offset f to a bank, and a bit a that is used to select the bank. [...]

Your problem is to determine the minimum running time of programs. In particular, given the number and size of the memory banks and a program to be executed, find the minimum number of instructions (which reference memory location and possibly set the BSR) that must be executed to run the program.

Input:

2 1
V1 V2 V1 V1 V2

Output:

6

Example for RE_PF

Question:

Read the following context: A statement after a two-hour emergency meeting at Stormont Castle... considered the Provisional IRA’s “bombing outrages”. [...]

Shootings in Belfast had continued last night. In one incident, a man was killed and five more were wounded... Later, a man was shot dead when he answered a knock at his door...

In the afternoon, during an hour of concentrated bombing at least 11 people died. They included two soldiers, a little girl and a messenger boy. [...] The worst explosion was at the busy Oxford Street bus station, where at least six people were killed.

Now please respond: Which time of attack caused more casualties, last night or afternoon?

- Procedure:**
1. What are the events containing victims or targets? (killed@228, died@354)
 2. What are the times in #1? (last night, afternoon)
 3. Which time of attack caused more people to die? (afternoon)

Answer:

afternoon

Example for IFE_CCL

Question:

Why are there 396 calories in a serving of 32 grams of unsalted almonds? Do not use commas in your response.

Answer:

```
[{'num_highlights': None, 'relation': None, 'num_words': None, 'num_placeholders': None, 'prompt_to_repeat': None, 'num_bullets': None, 'section_splitter': None, 'num_sections': None, 'capital_relation': None, 'capital_frequency': None, 'keywords': None, 'num_paragraphs': None, 'language': None, 'let_relation': None, 'letter': None, 'let_frequency': None, 'end_phrase': None, 'forbidden_words': None, 'keyword': None, 'frequency': None, 'num_sentences': None, 'postscript_marker': None, 'first_word': None, 'nth_paragraph': None}]
```

Example for IFE_EF

Question:

Create a JSON output that provides predictive analytics for the likelihood of success in a civil lawsuit... specifically regarding patent infringement... Please follow the structure as outlined below: [...]

Answer (Schema-Compliant JSON):

```
{
  "case_details": { "plaintiff": { "name": "TechCore Innovations", ... }, "defendant": { ... },
  "case_facts": { "patent_id": "US-9876543-B2", "technology_sector": "Semiconductors", "jurisdiction": "E.D. Texas" } },
  "historical_cases_analytics": { "total_cases_analyzed": 450, "win_rate_for_plaintiff": 0.42, "jurisdiction_success_rate": 0.55 },
  "predictive_outcome": { "success_likelihood": 0.68, "recommended_strategy": "Settle", "potential_damages_range": { "low": 5000000, "high": 12000000 } }
}
```

Example for KM_TK

Question:

How to evaluate Netflix’s \$82.7 billion acquisition of Warner Bros., HBO, and HBO Max? What impact will it have on the film and television industry?

Answer:

No_Info

Example for KM_NPK

Question:

What is the full transcript of the French representative’s remarks during a closed-door meeting of the UN Security Council discussing the Russia-Ukraine conflict?

Answer:

No_Info

Example for KE_IMC

Question:

The title “A3 motorway” is shared by multiple different works... Identify the precise genre... *Context:* thumb|290px|Glärner Alps... The A3 is a motorway in northeast Switzerland...

Answer:

Confoederatio Helvetica

Example for KE_FD

Question:

What is the real origin and earlier meaning of the word “gringo”?

Answer:

“Gringo” did not originate from song lyrics about “green grow”... It originally meant “foreigner”...

K.2 Prompts for Tasks

Prompt for RE_III

Role:

You are an expert Plain Language Specialist who rewrites complex texts to improve readability.

Task:

Analyze the provided source text without adding or removing factual details.

Output Format:

Output ONLY the simplified text string without explanations.

Example:**Input:**

Simplify the following text to improve its readability, ensuring its core meaning remains intact: “the land before time dvd the film explores issues of prejudice between the different species and the hardships they endure in their journey as they are guided by the spirit of littlefoot s mother.” Provide only the simplified text as the response.

Output:

The ‘Land Before Time’ DVD explores prejudice between species and the hardships they endure on their journey. They are guided by the spirit of Littlefoot’s mother.

Prompt for RE_PF

Role:

You are a precision information extraction engine who retrieves exact answers from knowledge graph triples.

Task:

Analyze the provided context facts and question to extract the exact answer entity.

Output Format:

Output ONLY the answer string without explanations.

Example:**Input:**

Read the following context: [‘Birbhum district located in the administrative territorial entity West Bengal, Birbhum district country India, Birbhum district located in the administrative territorial entity West Bengal Pradesh Congress Committee’] Now please respond: Which country is Birbhum district located in?

Output:

India

Prompt for KM_DSK

Role:

You are a professional medical health information query assistant.

Task:

You will receive a specific question regarding the public healthcare sector. Your task is to provide the correct answer. If the statement is entirely within the scope of your training data and you are certain of its accuracy, you must output: [TRUE]. If the statement is entirely within the scope of your training data and you are certain it is incorrect, you must output: [FALSE]. If the query involves any information not covered in your training data, or if you cannot verify its authenticity, you must output: [NO_INFO].

Output Format:

Return ONLY one of the following tokens: [TRUE], [FALSE], or [NO_INFO], with no explanations.

Example 1:**Input:**

Question: According to the latest 2025 research, can Alzheimer's disease be completely cured through novel gene therapy?

Output:

[NO_INFO]

Example 2:**Input:**

Question: "Can aspirin be used to prevent heart attacks?"

Output:

[TRUE]

Example 3:**Input:**

Question: "Does smoking cigarettes improve lung function?"

Output:

[FALSE]

Prompt for KM_TK

Role:

You are a meticulous data analyst specializing in numerical datasets.

Task:

You will receive a specific question about numerical data. Your task is to provide the correct answer. If the query involves real-world entities or facts you are confident about, please provide a brief and direct answer. If the query involves information not in your training data, you must output: [NO_INFO].

Output Format:

Return ONLY the answer or the token [NO_INFO], with no explanations.

Example 1:**Input:**

Question: As of the time I asked, what is China's current per capita GDP in US dollars?

Output:

[NO_INFO]

Example 2:**Input:**

Question: What was the number of applicants for the 2022 National College Entrance Examination?

Output:

11.93 million

Prompt for KE_EIC

Role:

You are a linguistic expert specialized in lexical semantics and word sense disambiguation.

Task:

Carefully analyze the contextual usage of the target word in both sentences. Decide whether the word conveys the same underlying sense in each context. Focus only on contextual interpretation without relying on external definitions or world knowledge.

Output Format:

Output a single word only: "Yes" or "No". Do not include any explanations.

Example 1:**Input:**

Question:

Sentence (A): After we leave the quarry, we intend to afforest the land and turn it into a nature reserve.

Sentence (B): Afforest the mountains.

Do both sentences use "afforest" with the same sense?

Output:

Yes

Example 2:**Input:**

Question:

Sentence (A): The surgeon closed the incision with a fine needle and thread.

Sentence (B): She threaded her way through the dense crowd at the market.

Do both sentences use "thread" with the same sense?

Output:

No

Prompt for KE_FD

Role:

You are an expert who specializes in judging whether there are errors in factual knowledge in various fields.

Task:

Judge whether the given question contains factual errors and provide explanations. Focus on examining different areas of knowledge to see if there are any misconceptions that need to be clarified or corrected.

Output Format:

Return ONLY the factual answer.

Example 1:**Input:**

Question: Did Coca-Cola invent the modern image of Santa Claus?

Output:

No. The red-suited Santa image existed before Coca-Cola used it in advertising.

Example 2:**Input:**

Question: How does linguistic history contradict the idea that "news" is an acronym for "North, East, West, South"?

Output:

News is not an acronym for compass directions. It developed in the 14th century as a special use of the plural for new (meaning new things), modeled after the French "nouvelles" or Latin "nova."

Prompt for IFE_EF

Role:

You are an expert in data serialization and strict text formatting.

Task:

Convert natural language instructions into precise formats. This encompasses both machine-readable data schemas (JSON, XML, YAML, CSV) and specific text structures (titles, bullet lists, highlights, headers, or section dividers).

Output Format:

Return ONLY the formatted string. Do not include any explanations.

Example:**Input:**

Generate a JSON object representing a book with the title “The Great Gatsby”, author “F. Scott Fitzgerald”, and year 1925. The keys must be “title”, “author”, and “year”.

Output:

```
{  
  "title": "The Great Gatsby",  
  "author": "F. Scott Fitzgerald",  
  "year": 1925  
}
```

Prompt for IFE_CCL

Role:

You are an expert at solving multi-constraint problems and adhere to strict formatting and content specifications.

Task:

Solve tasks that require deep reasoning and are subject to strict formatting rules. You must satisfy all structural constraints.

Output Format:

Return ONLY the final response. Do not include any explanations.

Example:**Input:**

Task: Find the product of 8 and 9.

Constraint Set:

1. Your response must be in all capital letters.
2. Do not use any numbers (digits) in the output; spell them out.
3. Do not use the letter ‘E’.

Output:

SIXTY-TWO

Prompt for Schema Normalizer

Role:

You are a data extraction agent specialized in converting unstructured text fragments into normalized Question-Answer schemas.

Task:

Analyze the input text to extract the core Question and Answer, stripping away all labels, instructional noise, and extra whitespace. Simultaneously, identify the specific Source (citation, title, or dataset name); if no source is explicitly mentioned, strictly set the value to null. If a question or answer is missing, set its respective field to null.

Output Format:

Return a single JSON object containing exactly three fields: Source, Question, and Answer. Do not include any explanations.

Example:**Input:**

“Wikipedia (Wiki-101): Plants convert sunlight into chemical energy via photosynthesis. Q: What is photosynthesis? A: the process of converting light energy into chemical energy.”

Output:

```
{"Source": "Wikipedia (Wiki-101)", "Question": "What is photosynthesis?", "Answer": "the process of converting light energy into chemical energy"}
```

Prompt for Evidence Retriever

Role:

You are an evidence verification agent responsible for validating whether each normalized QA instance is grounded in explicit textual evidence.

Task:

Given a normalized QA instance containing Source, Question, and Answer fields, use the provided Source to locate the corresponding content. Search for text segments that directly support or contradict the Answer. If found, extract the most relevant sentence or passage verbatim as Evidence. Do not paraphrase, infer missing facts, or use external information. If no such evidence is available, set the Evidence field to null.

Output Format:

Return a single JSON object with Evidence Source and Evidence. If no evidence is found, set the Evidence field to null.

Example:**Input:**

```
{"Source": "Wikipedia (Wiki-101)", "Question": "What is photosynthesis?", "Answer": "the process of converting light energy into chemical energy"}
```

Output:

```
{"Evidence Source": "Wikipedia (Wiki-101)", "Evidence": "Photosynthesis is the process by which plants convert light energy into chemical energy."}
```

Prompt for Type Classifier

Role:

You are a classification agent specialized in identifying the potential failure type that a question may induce when answered by a language model, based on a four-class taxonomy.

Task:

Inspect the question in the input QA pair and assign exactly one failure category that best characterizes the kind of mistake a language model is likely to make if it fails. Use the following four-category taxonomy based on the criteria defined in the evaluation guide: KE, KM, RE or IFE. Do not use multiple labels. Base your classification only on the question and its potential risk.

Output Format:

Return only the single category label (“KE”, “KM”, “RE”, or “IFE”) that you judge to be the most appropriate for the question. Do not provide any explanation or additional text.

Example:**Input:**

```
{ "Source": "CNN (2023-05)", "Question": "What causes hurricanes to rotate counterclockwise in the Northern Hemisphere?", "Answer": "Due to the Coriolis effect" }
```

Output:

```
{"Type": "KE" }
```

Prompt for Quality Scoring

Role:

You are a scoring agent designed to assess the quality of question-answer instances using three quality dimensions: Factuality, Discriminability, and Clarity.

Task:

Given an input QA triplet (Source, Question, Answer), assign a score from 1 to 10 for each of the following dimensions based on predefined evaluation criteria:

- **Factuality:** Evaluate whether the answer is accurate and well-grounded in the source, ensuring evidence consistency and no hallucinations.
- **Discriminability:** Evaluate whether the question cleanly targets a single failure type with no ambiguity or overlap across categories.
- **Clarity:** Evaluate whether the question and answer are linguistically clear, unambiguous, and well-phrased.

Each score must be an integer between 1 and 10, based on the criteria defined in the evaluation guide. Do not output explanations or justifications.

Output Format:

Return a single JSON object containing exactly three fields {Factuality, Discriminability, Clarity}.

Example:**Input:**

```
{  
  "Source": "CNN (2023-05)",  
  "Question": "What causes hurricanes to rotate counterclockwise in the Northern Hemisphere?",  
  "Answer": "Due to the Coriolis effect",  
  "Type": "KE"  
}
```

Output:

```
{10, 8, 10}
```