

UNDERSTANDING CONFORMAL FACTUALITY FOR RAG-BASED LLMs: NOVEL METRICS AND SYSTEMATIC INSIGHTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) are powerful generative models which can often generate responses that are plausible but not grounded in factual reality usually referred to as “hallucinations”. This is a challenge for using LLMs in applications that require answers that are factually correct. Two approaches have emerged as promising ways to mitigate this issue in the literature: (i) Conformal factuality filtering framework that provides statistical guarantee on the factual accuracy of claims in the final output, but cannot mitigate the hallucinations in the response generation and (ii) retrieval-augmented generation (RAG), which utilizes trusted knowledge bases as reference to guide the generation of response with the aim of reducing hallucinations but does not offer statistical guarantees. In this work, we unite these two approaches by integrating a conformal factuality framework with RAG and systematically study their performance to understand their strengths and limitations. We investigate the role of different key components: reference for generation and scoring functions, sensitivity to calibration data, LLM model capacity, reasoning and robustness to distractors. We propose three new metrics: *non-empty rate*, *non-vacuous empirical factuality*, and *sufficient correctness* to address limitations of standard factuality measures that fail to meaningfully capture usefulness of the output. Our experiments are comprehensive spanning three datasets (FactScore, MATH, and Natural Questions) and multiple model families and sizes. Our results show the importance of designing scoring functions and highlight trade-offs between correctness and informativeness that standard metrics fail to capture. Together, our findings provide insights that are practically useful and sheds light on the importance of re-thinking LLM factuality.

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities across open-domain question answering, reasoning, and scientific discovery (Brown et al., 2020; Guo et al., 2025; Zhang et al., 2025). Yet, a persistent barrier to their reliable deployment is the phenomenon of *hallucinations*: outputs that are fluent and confident but factually incorrect (Ji et al., 2023; Nadeau et al., 2024; Huang et al., 2025). Such errors are not merely cosmetic. In safety-critical settings, such as medicine, law, or finance, a single fabricated claim can erode trust, propagate misinformation, and incur high societal or financial costs. This makes hallucination mitigation one of the central challenges in advancing trustworthy LLMs.

A rich body of work has emerged to address this challenge with two main directions: (1) retrieval-augmented generation (RAG) and (2) conformal methods. RAG aims to prevent hallucinations by steering the generation process by grounding responses in trusted external knowledge sources. This is usually done by conditioning generation on retrieved passages (Lewis et al., 2020; Gao et al., 2023; Siriwardhana et al., 2023). While RAG reduces the likelihood of unsupported claims, it does not offer statistical guarantees on the factuality of the final response. Conformal prediction (CP) framework, on the other hand, aims to provide a statistical guarantee on the final output, which is often achieved via post-processing of the initial LLM response. CP frameworks usually first decompose the initial LLM output into atomic claims, score each with a given factuality scoring function, and filter those falling below a threshold determined using a calibration data (Mohri & Hashimoto,

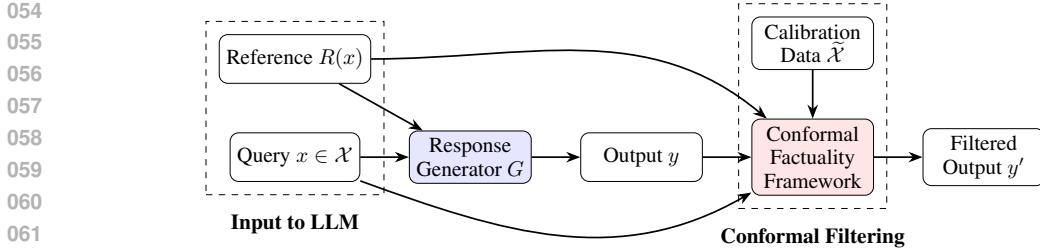


Figure 1: Overview of our framework. Given a query x and retrieved references $R(x)$, the Response Generator G produces an output y . The conformal factuality framework utilizes a separate calibration data to determine a threshold used to filter out information from the output y and yield y' .

2024; Cherian et al., 2024). This procedure provides formal coverage guarantees but often at the expense of informativeness, since aggressive filtering may yield empty or vacuous outputs. Furthermore, CP filtering cannot improve the usefulness or the accuracy of the LLM response, but only potentially remove hallucinations from the response.

Despite complementary strengths, these two paradigms have primarily been studied in isolation. It is natural to wonder whether combining them allows for significantly improving LLM factuality. There are a few recent works that combine RAG with CP methods (Li et al., 2023; Rouzrokh et al., 2024; Feng et al., 2025), but they are limited in their systematic analysis. A comprehensive understanding of the strengths and limitations of combining RAG and CP methods requires not only integrating the two frameworks but also developing principled ways to evaluate the performance. Standard metrics such as empirical factuality fail to capture the usefulness of the final answer, e.g., an empty answer is by default factually correct.

Our Contributions: We *systematically* investigate the integration of conformal factuality filtering with RAG-based LLMs to establish a unified framework for hallucination mitigation. Our comprehensive experiments span diverse datasets, model families, and scoring functions. Furthermore, we evaluate robustness against distribution shifts and distractors, shedding light on both the capabilities and limitations of this approach.

Beyond empirical analysis, we also address the limitations of current evaluation practices. To this end, we introduce three complementary metrics: *non-empty rate* and *non-vacuous empirical factuality*, which jointly capture both correctness and information retention, and *sufficient correctness*, which measures whether an output contains enough correct information to infer the final answer to the query. These novel metrics capture the trade-off between factual correctness and informativeness, providing practical insights and tools for future work on hallucination mitigation.

2 PROBLEM SETTING, DATASETS, MODELS, AND METRICS

Let $x \in \mathcal{X}$ denote an input query. Let $R(x)$ denote the reference material sufficient to answer x . We assume the existence of an oracle retriever that can retrieve $R(x)$ such that the true answer y^* is either in or can be deduced from $R(x)$. A response generator G , instantiated as a large language model (LLM), is prompted with both x and $R(x)$ to produce an output $y = G(x, R(x))$. Even with reference, LLMs can produce hallucinations in the generated response (Huang et al., 2025). The goal of the overall system is to create a final output y' such that each statement in y' is factually correct at a user’s expected level, e.g., 85%, while being useful in answering the query. Conformal filtering methods provide such statistical guarantees, which we describe below.

Conformal Filtering. Let $\tilde{\mathcal{X}}$ denote a calibration set disjoint from \mathcal{X} . For each query $x \in \tilde{\mathcal{X}}$, we obtain a set of claims $\{c_i\}$ by parsing the corresponding y using a parser P . These claims are then scored by a factuality scoring function f . Given a user-specified error level $\alpha \in (0, 1)$, the conformal filtering framework determines a threshold τ_α from the calibration set. Specifically, τ_α is chosen as the $\frac{[(n+1)(1-\alpha)]}{n}$ quantile of a set of candidate thresholds, where $n = |\tilde{\mathcal{X}}|$. For each calibration query x , the candidate threshold is defined as the smallest value of τ such that all claims surviving

above τ are factual: $\inf\{\tau \mid \forall c \in F(\tau), c \text{ is factual}\}$, where $F(\tau) = \{c \mid f(c) > \tau\}$ denotes the set of claims scoring above τ . The resulting guarantee is that $\mathbb{P}(\forall c \in F(\tau_\alpha), c \text{ is factual}) \geq 1 - \alpha$. At inference time, each generated output y is decomposed into atomic claims by a parser P , yielding $C(y) = \{c_i\}_{i=1}^k$. Each claim c_i is scored with f , and those exceeding the threshold are retained: $C'(y) = \{c_i \mid f(c_i) > \tau_\alpha, i = 1, \dots, k\}$. Finally, the retained claims $C'(y)$ are merged by a merger M into a single filtered response: $y' = M(C'(y))$. Figure A.1.3 in the Appendix illustrates the whole pipeline of the conformal filtering system.

2.1 SCORING FUNCTIONS

Scoring functions f are the core component of the conformal factuality framework, as they determine whether a claim is retained or filtered. We study two main families of scoring functions:

(a) Entailment-based scorers are natural language inference (NLI) models to assess whether the reference text supports a claim. We use both document-level and sentence-level variants. For **document-level entailment**, the entailment score is computed directly between the entire reference $R(x)$ and the claim c_i . We use an entailment model (Laurer et al., 2022) trained on the DocNLI dataset (Yin et al., 2021) for this.

In the sentence-level setting, entailment scores are computed between the claim c_i and each sentence in $R(x)$ and then aggregated in two different ways. The first is **conservative entailment**, where a claim is marked as contradictory if any sentence in the reference contradicts it. It is marked as entailed only if there are no contradictions and there is at least one sentence supporting it. The second is **average entailment**, which averages the scores of all non-neutral sentence-level comparisons. For sentence-level entailment, we use `roberta-large-mnli` as the entailment model (Liu et al., 2019).

(b) LLM-based scorers. In this family, we prompt a language model¹ to assign a factuality score to each claim. We explore the design space by varying five dimensions: (i) the inclusion of retrieved references $R(x)$, (ii) evidence highlighting within $R(x)$, (iii) the use of Chain-of-Thought (CoT) reasoning, (iv) output granularity (continuous $[0, 1]$ vs. Boolean), and (v) evaluation consistency (single generation vs. averaging over five independent generations).

2.2 DATASETS

We perform evaluations on three datasets that span open-ended summarization, mathematical reasoning, and question answering tasks. This diversity allows us to assess both the factual reliability and the task-level utility of our approach.

- **FactScore** dataset (Min et al., 2023) consists of 601 individuals, each paired with a Wikipedia page. Queries are: “Tell me a paragraph biography about [person],” where the reference $R(x)$ is the Wikipedia page of the person. Because no canonical ground-truth answers are provided, this dataset is well-suited for evaluating factuality in *open-ended generation* and for testing whether models can produce faithful summaries grounded in external references. We additionally consider **FactScore Rare**, a subset of 198 queries focusing on less well-known individuals. This subset probes the robustness of models when the model’s parametric knowledge is not enough to answer the question, a regime where hallucinations are more likely, if a reference is not given.
- **MATH** dataset (Hendrycks et al., 2021) contains 12,446 competition-style mathematics problems spanning five difficulty levels and seven categories. Each problem provides a question x and a ground-truth answer y^* . To construct reference materials, we prompted `gpt-5-nano` to generate prerequisite knowledge relevant to solving each problem, which serves as $R(x)$ ².
- **Natural Questions (NQ)** dataset (Kwiatkowski et al., 2019) consists of 10K real-world queries collected from search engines. Each query is annotated with both a long answer and a short answer; we use the long answer as the reference $R(x)$ and the short answer as the ground-truth.

Together, these datasets provide coverage over distinct capabilities: factual summarization (FactScore), mathematical reasoning (MATH), and reference-based question answering (NQ).

¹Prompts provided in Appendix A.5.8.

²Prompts provided in Appendix A.5.14

2.3 LANGUAGE MODELS

We evaluate our framework across both open-source and proprietary language models in order to systematically evaluate parts and whole of the Factuality and RAG pipeline under varying architectures, reasoning modes, and parameter scales.

Open-source models. Our open-source suite includes multiple families. The Qwen3 models (Yang et al., 2025) (0.6B, 4B, and 8B parameters) are evaluated both in their base form and in a reasoning-enabled variant, Qwen3-Think, where reasoning with the `<think></think>` tag is enabled. This contrast allows us to study whether reasoning-oriented training improves factuality scoring and filtering. To broaden architectural diversity, we also include Llama-3.2-Instruct (Dubey et al., 2024) (1B and 3B parameters), SmoLLM2-Instruct (Allal et al., 2025) (135M, 360M, and 1.7B parameters), and gpt-oss-20b³ (Agarwal et al., 2025) (21B parameters with 3.6B active).

Proprietary model. In addition to open-source models, we employ gpt-5-nano (gpt, 2025), a proprietary model, primarily as an auxiliary component. We use it for tasks such as claim parsing, factuality labeling, and claim merging, where higher accuracy in controlled subtasks helps construct reliable evaluation pipelines. Importantly, the proprietary model is not the subject of our evaluation, but rather supports the analysis of the open-source models.

The chosen models are designed to probe three orthogonal dimensions: (i) *model architecture*, by comparing across families; (ii) *reasoning capability*, by contrasting Qwen3 with Qwen3-Think; and (iii) *model scale*. This diversity allows us to assess how each factor influences factuality scoring, and to highlight regimes where smaller, more efficient models suffice.

2.4 EVALUATION METRICS

We evaluate the performance using both established factuality measures and new metrics designed to capture aspects of informativeness that existing metrics overlook. We use three widely used criteria:

- **Empirical Factuality (EF)** measures the fraction of outputs y' in which all retained claims $C'(y)$ are factual. By convention, an empty claim set $C'(y) = \emptyset$ is treated as factual, which can artificially inflate EF when filtering is overly aggressive. *Higher EF is better*, as it indicates stronger factual reliability.
- **Power** quantifies the average proportion of true claims retained. *Higher Power is better*, since it means fewer correct claims are lost.
- **False Positive Rate (FPR)** measures the fraction of non-factual claims that survive. *Lower FPR is better*, as it reflects stronger suppression of hallucinations.

These metrics only provide limited insight into the overall usefulness of the final answer. Each of the statements in an LLM response could be factually correct while still not being informative enough to answer the input query. Furthermore, a vacuous or empty output is factually correct by definition, but is not useful, and when the reference does contain information to provide a correct answer, an empty answer is an indication of failure. While EF, Power, and FPR capture factuality and error rates, they fail to penalize vacuous but “factual” outputs. To address this, we propose the following novel evaluation metrics.

- **Non-empty Rate (NR)**, the fraction of outputs that preserve at least one claim. *Higher NR is better*, rewarding informative responses rather than empty ones.
- **Non-vacuous Empirical Factuality (NvEF)**, which computes EF only over non-empty outputs. *Higher NvEF is better*, reflecting factuality conditional on informativeness.
- **Correctness** measures the fraction of outputs equivalent to the ground-truth answer y^* . *Higher Correctness is better*, though this metric is intentionally strict and more applicable to the dataset with ground truth answers.
- **Sufficient Correctness (SC)**, which holds if the filtered output contains enough correct information, relative to a reference $R(x)$, to infer the correct answer. *Higher SC is better*, as it captures end-task utility. This metric is inspired by sufficient context introduced in (Joren et al., 2025).

NR, NvEF are claim-level, while correctness and sufficient correctness are at the final task-level outcome. Together, these metrics balance factual reliability, informativeness, and task-level utility, ensuring that evaluation reflects not only safety (removing hallucinations) but also usefulness (retaining an adequate signal for the end task).

³Results are deferred to Appendix A.2.8

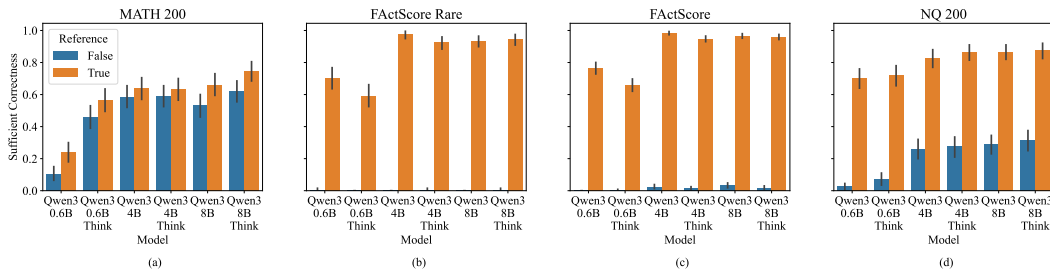


Figure 2: Sufficient correctness (SC) of Qwen3 models (0.6B, 4B, 8B) on four datasets (MATH-200, FActScore-Rare, FActScore, NQ-200), with and without access to references. Across model sizes and datasets, providing references consistently improves generation quality.

We now design a series of experiments to systematically analyze: (i) the role of references, (ii) the design of scoring functions, and (iii) robustness to distributional shifts and adversarial inputs.

3 IMPACT OF REFERENCES

We begin by isolating the role of references, asking: *How much does retrieved reference improve generation quality before filtering is even applied?* To study this, we evaluate outputs y produced by the response generator G under two conditions: query-only generation $y = G(x)$ and query-plus-reference generation $y = G(x, R(x))$. We then measure sufficient correctness and compare across FActScore, FActScore-Rare, **MATH-200** (a 200-example subset of MATH), and **NQ-200** (a 200-example subset of Natural Questions). These datasets help provide insights into whether conditioning on $R(x)$ improves different types of tasks: factual summarization, reasoning, and question answering.

Results. Figure 2 compares sufficient correctness (SC) with and without references for the Qwen3 family. Note that the FActScore and NQ datasets, reference plays a huge role, which highlights its importance when the model is not able to have all information memorized. On MATH, the gap between the two settings narrows as model size increases, while overall SC improves, suggesting that larger models possess stronger reasoning ability and can better leverage reference material. For Qwen3-0.6B, enabling *think* reasoning produces a large jump in SC, further highlighting the role of reasoning capacity. In contrast, for FActScore and NQ, where math reasoning plays a more minor role, enabling *think* has little effect. Lastly, we observe that there are diminishing returns in terms of increasing the model size in the Qwen3 model family. The improvement in terms of SC when we increase the model size from 4B to 8B is smaller compared to when we improve the model size from 0.6B to 4B. Overall, providing references consistently improves generation quality. We observe similar trends for SmolLM2 and Llama-3.2 (see Figure 12 and 13 in the Appendix), and therefore provide references to the generator in all subsequent experiments.

4 SCORING FUNCTIONS

Having shown that reference improves baseline quality, we now turn to conformal filtering, which offers statistical guarantees. Its backbone is the scoring function for factuality. In this section, we examine how prompting strategies, references, model choice, and different scores affect performance.

4.1 PROMPTING STRATEGIES FOR LLM SCORERS

We begin by examining the effect of prompting variations for LLM-based scoring functions, including: (i) highlighting supporting evidence, (ii) enabling chain-of-thought reasoning, (iii) scalar vs Boolean scoring, and (iv) consistency averaging. References are always provided to both f and G . We run this experiment on **FActScore**, **MATH-1K** (a 1,050-example subset of MATH), and **NQ-1K** (a 1,000-example subset of Natural Questions).

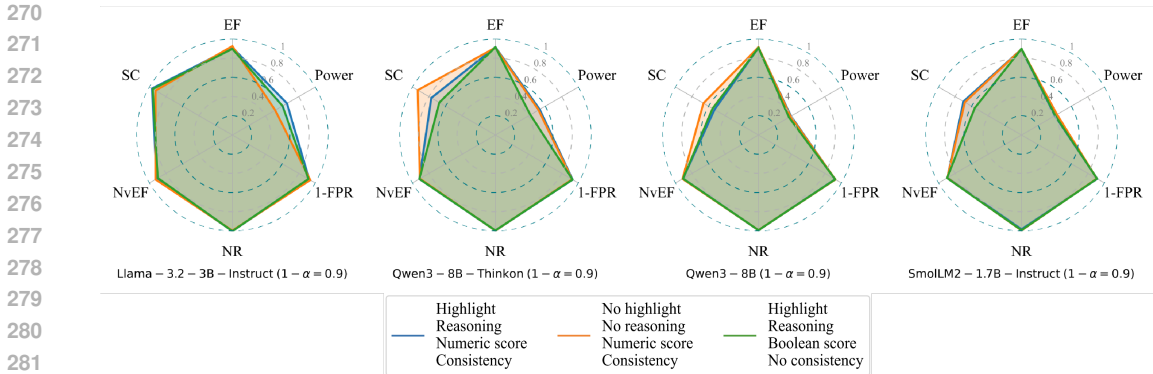


Figure 3: We evaluated 16 prompting strategies across three LLMs, presenting three representative approaches for clarity. Results demonstrate that: (i) prompting models to generate numeric scores consistently outperforms boolean scoring; (ii) sampling multiple responses uniformly improves performance; however, (iii) incorporating chain-of-reasoning or response highlighting as grounding mechanisms does not yield reliable performance gains across models.

Results. Figure 3 summarizes the effect of prompting strategies on model confidence scores across three LLM families. While no single strategy emerges as universally optimal, several consistent patterns are observed. First, instructing models to output numeric scores reliably outperforms boolean judgments, suggesting that scalar responses better capture gradations of uncertainty. Second, sampling multiple responses provides consistent gains, indicating that aggregation reduces variance and stabilizes confidence estimates. In contrast, more elaborate mechanisms, such as chain-of-reasoning prompts or response highlighting, fail to deliver consistent improvements and sometimes even degrade performance. Our findings demonstrate that generic prompting heuristics prove inadequate for conformal factuality tasks, necessitating architecture-specific prompting strategies and highlighting the importance of principled, fine-grained optimization approaches such as TextGrad (Yuksekonul et al., 2024) for future research.

4.2 ROLE OF REFERENCES IN SCORING

In this section, we evaluate how providing references to the LLM-based scoring function improves it. Outputs $y = G(x, R(x))$ are generated using gpt-5-nano, while open-source models serve as scorers. This is because API calls are faster than using our GPU for inference. When varying reference access, we fix all other prompt settings (CoT enabled, scalar scoring, and consistency averaging). Experiments are conducted on FActScore, MATH-1K, and NQ-1K.

Results. Figure 4 shows the performance under various metrics for the model confidence scores with and without reference provided to the scoring function for the MATH-1K dataset using Qwen3-4B as the scorer. We observe that when a reference is introduced to the scoring function, the power is consistently higher than in the case where no reference is given to the scoring function. Similarly, for the non-empty rate (NR), the scoring functions with reference have an advantage when the target factuality is of a larger size. These results show the benefit of feeding the reference to the scoring function. Therefore, in the subsequent experiments, we provide a reference to our scoring function.

4.3 MODEL CHOICE FOR LLM-BASED SCORERS

We compare different open-source models as scoring functions f , using gpt-5-nano as the generator G . All prompts include references, require highlighting and CoT, produce scalar scores, and apply consistency averaging. This experiment assesses how the scorer model family and scale affect factuality filtering.

Results. Our experiments reveal that scaling LLMs does not guarantee improved confidence calibration in conformal factuality, as shown in Figure 5. While the Llama series exhibits a consistent gain on model confidence score with larger models, this trend breaks for other families. The SmolLM2 series shows no systematic benefit, and in the case of Qwen3, scaling even degrades

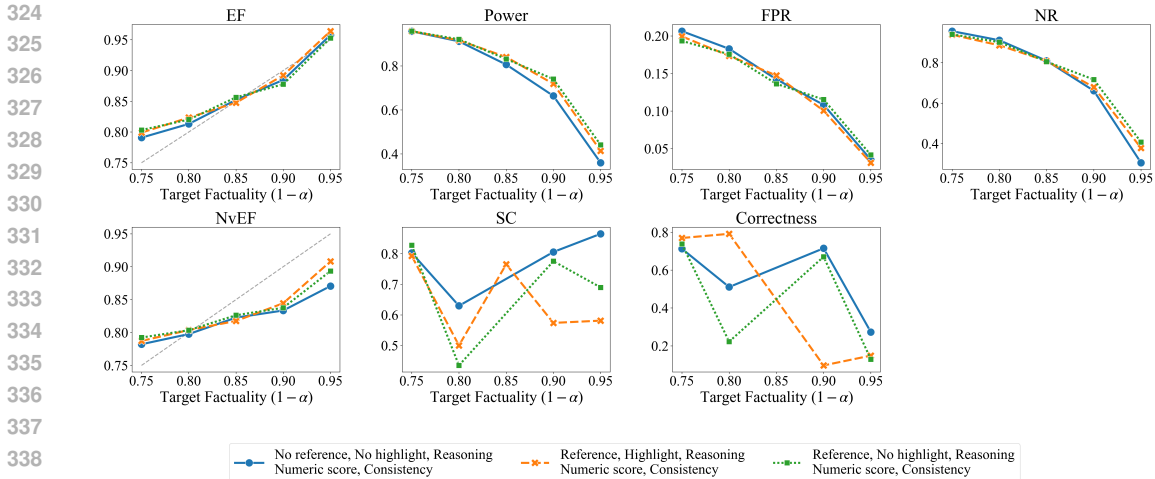


Figure 4: Performance of model confidence score on MATH-1K dataset with and without reference provided to scoring functions using Qwen3-4B as the LLM-based scorer.

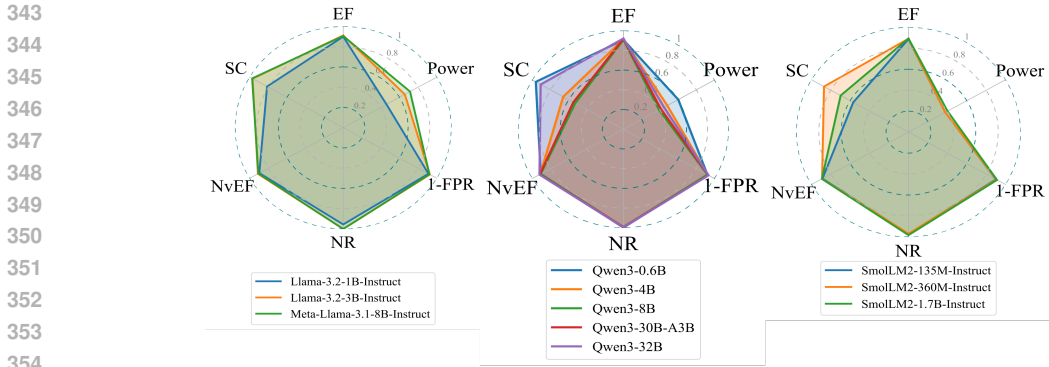


Figure 5: We conducted experiments across three model series of varying scales. Results demonstrate that increasing model size does not consistently improve performance in conformal factuality. While scaling the Llama series yields consistent improvements in confidence score accuracy, this pattern does not hold for the Qwen3 and SmoLLM2 series. Notably, for the Qwen3 series, scaling may degrade performance.

performance. These heterogeneous scaling behaviors challenge the assumption that model size universally improves calibration quality, suggesting that progress in conformal factuality will require more fine-grained design beyond naive scaling.

4.4 DIFFERENT FAMILY OF SCORING FUNCTIONS

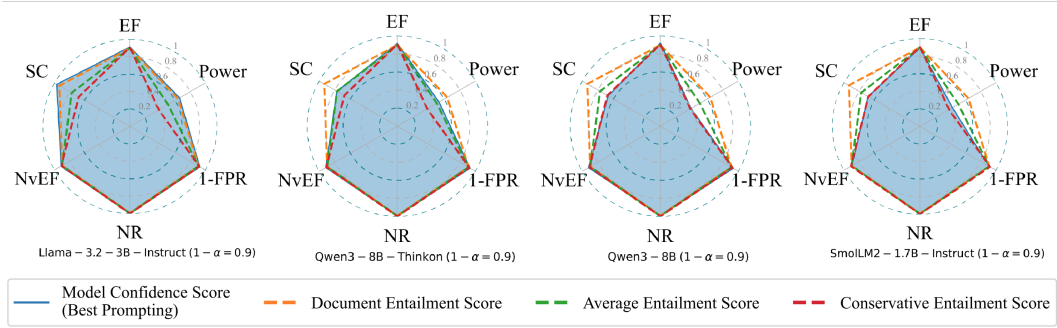
Beyond model confidence scores, we also examine entailment-based scoring functions. We compare confidence-based filtering with entailment-based filtering on FActScore, MATH-1K, and NQ-1K datasets to assess whether entailment signals offer additional advantages in factuality filtering.

Results. Figure 6 presents a direct comparison between model confidence scores⁴ and entailment-based scoring. Notably, entailment-based scores consistently match or exceed model confidence, despite the entailment model being substantially smaller than the target LLMs⁵. Together with our scaling analysis 4.3, these findings challenge the assumption that larger models inherently yield

⁴Under the best prompting strategy we identified in Section 4.1 for each LLM.

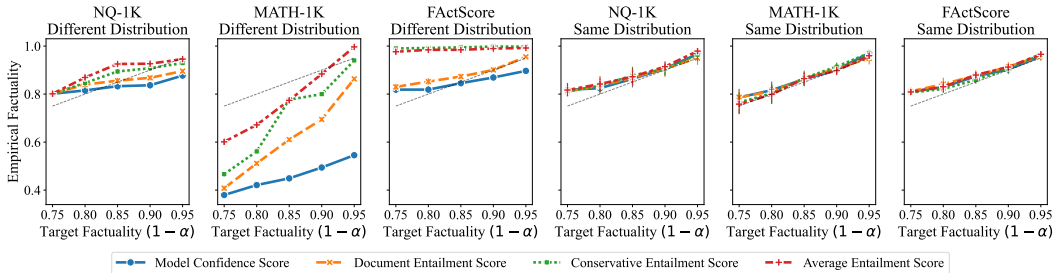
⁵We employ DeBERTa and RoBERTa to compute entailment scores, achieving computational efficiency gains of more than two orders of magnitude compared to LLM-based confidence scoring methods. For detailed comparisons of model parameters and computational complexity, see Table 1.

378
379
380
381
382
383
384
385
386
387
388
389



390
391
392
393
394
395
396
397
398
399
400
401
402
403

Figure 6: Using the best-performing prompt for each LLM identified in Section 4.1, we compare entailment-based scores against the model confidence score. Across our evaluated settings, the entailment-based scores match or exceed the model-confidence baseline, despite the entailment model being substantially smaller than the target LLMs.



404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420

Figure 7: Empirical factuality (EF) on FActScore, MATH-1K, and NQ-1K, under two calibration settings: (i) calibration claims generated in Mohri & Hashimoto (2024), which comes from a different distribution (ii) calibration claims drawn from the same distribution as the test data. We use Qwen3-4B as the scoring function. The results show how the distribution shift in the calibration set affects conformal factuality guarantees.

421
422
423
424
425
426
427
428
429
430
431

better-calibrated confidence estimates, instead, highlighting that targeted, lightweight verifiers can deliver both computational efficiency and superior performance.

5 ROBUSTNESS OF SCORING FUNCTIONS

Although scoring functions differ in their design and accuracy, their real-world usefulness depends on robustness: can they maintain reliability under distribution shifts or adversarial perturbations? We therefore stress-test our framework under challenging conditions.

5.1 ROBUSTNESS TO CALIBRATION DISTRIBUTION SHIFT

Conformal factuality assumes calibration and test data are exchangeable. We evaluate robustness when this assumption is violated by using mismatched calibration data. Specifically, we compare calibration on (i) 50 human-annotated claims from Mohri & Hashimoto (2024) (gpt-4 generated claims, denoted as MH) versus (ii) 50 randomly selected queries from the held-out test half. We evaluate empirical factuality using Qwen3-4B, SmolLM2-360M-Instruct, and Llama-3.2-3B-Instruct across FActScore, MATH-1K, and NQ-1K datasets.

Results. Figure 7 compares empirical factuality (EF) under different sources of calibration data, with test outputs generated and scored by Qwen3-4B. When calibration data come from a different distribution (generated by gpt-4), EF sometimes falls below the target level (grey line) for certain datasets and scoring functions. While the entailment based scorers are robust to this distribution shift, this robustness is not consistent across language models: switching G and f to SmolLM2-360M-Instruct or Llama-3.2-3B-Instruct yields different behaviors for the

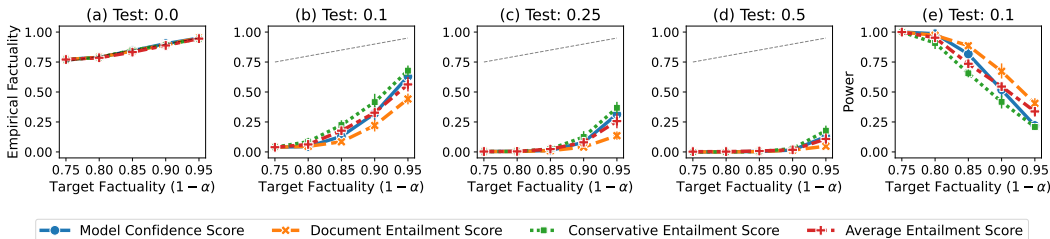


Figure 8: (a) - (d) Empirical factuality versus target factuality ($1 - \alpha$) for Qwen3-4B on the FActScore dataset. Each panel corresponds to a different distractor injection rate in the test set (0.0, 0.1, 0.25, 0.5). The gray dashed diagonal represents perfect calibration ($y = x$), where empirical factuality matches the target. (e) Power versus target factuality ($1 - \alpha$) when the injection rate is 0.1.

scoring functions that seem robust in the Qwen3-4B setting (see Figures 29 and 28 in the appendix). This shows that the CP filtering is sensitive to the distribution shifts between the claims during test time and the calibration set. Note that this arises even when the LLM used to generate the response is switched for the same dataset. So, the calibration data should be collected using the same LLM that the system is going to be deployed with.

5.2 ROBUSTNESS TO DISTRACTORS

In real-world scenarios, the LLM outputs may contain claims that are plausible but incorrect. This may be due to the fact that LLMs themselves are easily distracted by irrelevant information (Shi et al., 2023). To simulate this setting, we replace a proportion of factual claims in each test query with distractors generated by gpt-5-nano. We then evaluate whether the conformal factuality framework can reliably distinguish correct content from these plausible but non-factual claims. Experiments are conducted on FActScore, MATH-1K, and NQ-1K datasets. We defer the discussion of how we create the distractors to Appendix A.3.

Results. Figure 8 (a)-(d) show that as the distractor rate increases, empirical factuality (EF) drops sharply. This degradation occurs because adding distractors to the test set likely violates the exchangeability assumption underlying the conformal factuality framework, causing EF to fall below the target level. Although EF can increase when the target factuality is set very high, this comes at the cost of a substantial loss in power, as shown in Figure 8 (e). Our results indicate that CP filtering with current scoring functions is not robust to distractors, underscoring the need for improved scoring functions that maintain robustness in their presence.

5.2.1 CAN DISTRACTION-AWARE CALIBRATION HELP?

A natural way to mitigate the effect of distractors is to anticipate their presence by using distraction-aware calibration. To study the efficacy of this approach, we extend the previous setting by introducing distractors not only into the test set but also into the calibration set. This models potential distribution shifts caused by distractor content and evaluates whether conformal filtering remains reliable when calibration data include such claims. Since the true level of distractors in practice is unknown, we vary the amount of distractors in calibration and test sets independently. This setup enables us to assess both under-estimation and over-estimation of distractor prevalence. Experiments are conducted on FActScore, MATH-1K, and NQ-1K datasets.

Results. In Figure 9, we show the result for the setting with a test set with a distractor proportion of 0.25 and varying the levels of distractors in the calibration set. When the fraction of distractors is underestimated, the EF is still far below the target EF (Figure 9 (a)). From Figure 9 (b)-(c), we can see that introducing a large enough fraction of distractors to the calibration set can bring up the empirical factuality. However, we note that this incurs a high cost on the non-empty rate. As we see in Figure 9 (d)-(e), the non-empty rate drops significantly when distractors are introduced to the calibration set. When both calibration and test contain no distractors, we have a much higher non-empty rate compared to the case where we inject 25% distractors into both the calibration and test sets. This happens likely due to the thresholds found by the CP framework becoming more stringent.

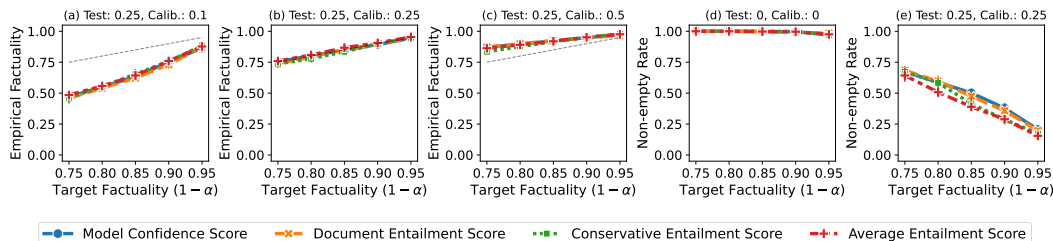


Figure 9: (a) - (c) Empirical factuality versus target factuality ($1-\alpha$) for Qwen3-4B on the FActScore dataset when the test set is injected with 25% distractors. We vary the proportion of distractors in the calibration set from (0.1, 0.25, 0.5). As the proportion matches, we can see that the empirical factuality rises to the $y = x$ line. (d)-(e) Comparison of non-empty rate when the test and calibration sets contain no distractors and contain 25% of the distractors. Although introducing distractors to the calibration set can achieve target factuality, the non-empty rate suffers.

Therefore, the scoring functions cannot distinguish the distractors that are factually incorrect from the factually correct claims.

6 RELATED WORKS

Many studies show that LLMs are prone to hallucination Nadeau et al. (2024); Huang et al. (2025) despite their impressive capabilities in summarization, dialogue, and coding Achiam et al. (2023); Zhang et al. (2024); Nam et al. (2024).

Retrieval-augmented generation (RAG) improves LLM performance on knowledge-intensive tasks by grounding responses in retrieved external context Lewis et al. (2020); Joren et al. (2025); Gao et al. (2023). While RAG improves the LLM generation, the output can still contain hallucinations Huang et al. (2025). Moreover, LLM may not utilize the reference well, especially that information in the middle, a phenomenon known as lost-in-the-middle Liu et al. (2023); Ravaut et al. (2023); Chen et al. (2023); Tang et al. (2023), which may be caused by the use of rotary positional embedding (RoPE) Su et al. (2024); Huang et al. (2025). Another reason for this to happen is that in the pretraining data, the most salient information is located at the beginning or the end, rather than in the middle Ravaut et al. (2023); Huang et al. (2025).

Conformal prediction (CP) to filter non-factual content from LLM outputs has emerged as a promising strategy Vovk et al. (2005); Angelopoulos & Bates (2021); Mohri & Hashimoto (2024); Cherian et al. (2024). These methods not only improve factuality but also offer a statistical guarantee: $\mathbb{P}(\text{Output is factual}) \geq 1 - \alpha$. For instance, Mohri & Hashimoto (2024) introduced *conformal factuality*, which scores claims in the output and removes those below a threshold.

7 CONCLUSION

We systematically investigated the role of key components in hallucination mitigation, including the use of references for generation and scoring, sensitivity to calibration data, model capacity, reasoning ability, and robustness to distractors. To address limitations of standard factuality measures, we proposed three new metrics: *non-empty rate*, *non-vacuous empirical factuality*, and *sufficient correctness*, which better capture the usefulness of filtered outputs. Our experiments span three diverse datasets (FActScore, MATH, and Natural Questions) and three model families. Our results highlight the importance of scoring function design and reveal trade-offs between correctness and informativeness that standard metrics overlook. Overall, our findings provide practical insights and underscore the need to rethink how LLM factuality is measured and enforced.

REFERENCES

Aug 2025. URL <https://openai.com/index/introducing-gpt-5/>.

- 540 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
541 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
542 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 543
- 544 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K
545 Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv*
546 *preprint arXiv:2508.10925*, 2025.
- 547 Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo,
548 Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav,
549 et al. Smollm2: When smol goes big—data-centric training of a small language model. *arXiv*
550 *preprint arXiv:2502.02737*, 2025.
- 551
- 552 Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and
553 distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- 554
- 555 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
556 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
557 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 558 Hung-Ting Chen, Fangyuan Xu, Shane Arora, and Eunsol Choi. Understanding retrieval augmenta-
559 tion for long-form question answering. *arXiv preprint arXiv:2310.12150*, 2023.
- 560
- 561 John Cherian, Isaac Gibbs, and Emmanuel Candes. Large language model validity via enhanced
562 conformal prediction methods. *Advances in Neural Information Processing Systems*, 37:114812–
563 114842, 2024.
- 564
- 565 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,
566 2025. URL <https://arxiv.org/abs/2501.12948>.
- 567
- 568 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
569 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
570 *arXiv e-prints*, pp. arXiv–2407, 2024.
- 571
- 572 Naihe Feng, Yi Sui, Shiyi Hou, Jesse C Cresswell, and Ga Wu. Response quality assessment for
573 retrieval-augmented generation via conditional conformal factuality. In *Proceedings of the 48th*
574 *International ACM SIGIR Conference on Research and Development in Information Retrieval*,
575 pp. 2832–2836, 2025.
- 576
- 577 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,
578 Haofen Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A
579 survey. *arXiv preprint arXiv:2312.10997*, 2(1), 2023.
- 580
- 581 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu
582 Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforc-
583 e-ment learning. *Nature*, 645(8081):633–638, 2025.
- 584
- 585 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced bert
586 with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- 587
- 588 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
589 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
590 *preprint arXiv:2103.03874*, 2021.
- 591
- 592 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
593 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language
594 models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information*
595 *Systems*, 43(2):1–55, 2025.
- 596
- 597 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
598 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM*
599 *computing surveys*, 55(12):1–38, 2023.

- 594 Hailey Joren, Jianyi Zhang, Chun-Sung Ferng, Da-Cheng Juan, Ankur Taly, and Cyrus Rashtchian.
595 Sufficient context: A new lens on retrieval augmented generation systems. In *The Thirteenth*
596 *International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Jjr2Odj8DJ>.
597
- 598 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
599 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a
600 benchmark for question answering research. *Transactions of the Association for Computational*
601 *Linguistics*, 7:453–466, 2019.
602
- 603 Moritz Laurer, Wouter van Atteveldt, Andreu Salleras Casas, and Kasper Welbers. Less annotating,
604 more classifying – addressing the data scarcity issue of supervised machine learning with deep
605 transfer learning and BERT-NLI, June 2022. URL <https://osf.io/74b8k>. Preprint, Open
606 Science Framework.
- 607 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
608 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented gener-
609 ation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:
610 9459–9474, 2020.
- 611 Shuo Li, Sangdon Park, Insup Lee, and Osbert Bastani. Traq: Trustworthy retrieval augmented
612 question answering via conformal prediction. *arXiv preprint arXiv:2307.04642*, 2023.
613
- 614 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
615 Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint*
616 *arXiv:2307.03172*, 2023.
- 617 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
618 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
619 approach. *arXiv preprint arXiv:1907.11692*, 2019.
620
- 621 Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,
622 Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual
623 precision in long form text generation. *arXiv preprint arXiv:2305.14251*, 2023.
- 624 Christopher Mohri and Tatsunori Hashimoto. Language models with conformal factuality guar-
625 antees. In *Proceedings of the 41st International Conference on Machine Learning*, ICML’24.
626 JMLR.org, 2024.
- 627 David Nadeau, Mike Kroutikov, Karen McNeil, and Simon Baribeau. Benchmarking llama2, mis-
628 tral, gemma and gpt for factuality, toxicity, bias and propensity for hallucinations. *arXiv preprint*
629 *arXiv:2404.09785*, 2024.
630
- 631 Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. Using
632 an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International*
633 *Conference on Software Engineering*, pp. 1–13, 2024.
- 634 Mathieu Ravaut, Aixin Sun, Nancy F Chen, and Shafiq Joty. On context utilization in summarization
635 with large language models. *arXiv preprint arXiv:2310.10570*, 2023.
636
- 637 Pouria Rouzrokh, Shahriar Faghani, Cooper U Gamble, Moein Shariatnia, and Bradley J Erickson.
638 Conflare: conformal large language model retrieval. *arXiv preprint arXiv:2404.04287*, 2024.
- 639 Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael
640 Schärli, and Denny Zhou. Large language models can be easily distracted by irrelevant context.
641 In *International Conference on Machine Learning*, pp. 31210–31227. PMLR, 2023.
642
- 643 Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and
644 Suranga Nanayakkara. Improving the domain adaptation of retrieval augmented generation (rag)
645 models for open domain question answering. *Transactions of the Association for Computational*
646 *Linguistics*, 11:1–17, 2023.
- 647 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-
hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

648 Raphael Tang, Xinyu Zhang, Xueguang Ma, Jimmy Lin, and Ferhan Ture. Found in the middle:
649 Permutation self-consistency improves listwise ranking in large language models. *arXiv preprint*
650 *arXiv:2310.07712*, 2023.

651 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
652 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly
653 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

654 Vladimir Vovk, Alexander Gammernan, and Glenn Shafer. *Algorithmic learning in a random world*,
655 volume 29. Springer, 2005.

656 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
657 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
658 *arXiv:2505.09388*, 2025.

659 Wenpeng Yin, Dragomir Radev, and Caiming Xiong. Docnli: A large-scale dataset for document-
660 level natural language inference. *arXiv preprint arXiv:2106.09449*, 2021.

661 Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and
662 James Zou. Textgrad: Automatic "differentiation" via text, 2024. URL <https://arxiv.org/abs/2406.07496>.

663 Yanbo Zhang, Sumeer A Khan, Adnan Mahmud, Huck Yang, Alexander Lavin, Michael Levin,
664 Jeremy Frey, Jared Dunnmon, James Evans, Alan Bundy, et al. Exploring the role of large lan-
665 guage models in the scientific method: from hypothesis to discovery. *npj Artificial Intelligence*, 1
666 (1):14, 2025.

667 Yang Zhang, Hanlei Jin, Dan Meng, Jun Wang, and Jinghua Tan. A comprehensive survey on
668 process-oriented automatic text summarization with exploration of llm-based methods. *arXiv*
669 *preprint arXiv:2403.02901*, 2024.

670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A APPENDIX

A.1 EXTENDED RELATED WORKS

A.1.1 RETRIEVAL-AUGMENTED GENERATION

Retrieval-augmented generation (RAG) improves LLM performance on knowledge-intensive tasks by grounding responses in retrieved external context Lewis et al. (2020); Joren et al. (2025); Gao et al. (2023). Formally, given a query x , a retriever R returns context $R(x)$ that supplements the model’s parametric knowledge, guiding generation toward more factually correct outputs. In this work, we assume access to an oracle retriever that always provides relevant, accurate references. While RAG is powerful, it does not provide statistical guarantees on the factuality of its outputs.

A.1.2 CONFORMAL PREDICTION

Despite their impressive capabilities in summarization, dialogue, and coding Achiam et al. (2023); Zhang et al. (2024); Nam et al. (2024), LLMs are prone to hallucination Nadeau et al. (2024); Huang et al. (2025). One promising mitigation strategy applies conformal prediction to filter non-factual content from LLM outputs Vovk et al. (2005); Angelopoulos & Bates (2021); Mohri & Hashimoto (2024); Cherian et al. (2024). These methods not only improve factuality but also offer a statistical guarantee: $\mathbb{P}(\text{Output is factual}) \geq 1 - \alpha$. For instance, Mohri & Hashimoto (2024) introduced *conformal factuality*, which scores claims in the output and removes those below a threshold calibrated on held-out data. The choice of scoring function f is therefore critical to the effectiveness of the framework.

A.1.3 CONFORMAL PREDICTION AND RAG

Recent work has begun integrating conformal prediction into retrieval-augmented generation (RAG) to provide statistical reliability guarantees. TRAQ Li et al. (2023) applies conformal prediction at both the retriever and generator stages, ensuring with high probability that a semantically correct answer is included in the output set. Conformal-RAG Feng et al. (2025) instead operates at the sub-claim level, filtering unreliable statements to guarantee factuality across domains. Conflare Rouzrokh et al. (2024) focuses on the retrieval stage, calibrating similarity thresholds so that retrieved contexts contain the true answer with user-specified confidence. While these approaches provide important coverage guarantees at different stages of the pipeline, they largely assess correctness in isolation. In contrast, our work unifies retrieval and filtering and introduces new metrics—non-empty rate, non-vacuous empirical factuality, and sufficient correctness—that explicitly capture the trade-off between correctness and informativeness, which existing frameworks do not address.

A.2 EXTENDED RESULTS

A.2.1 IMPACT OF REFERENCES

Section 3 in the main paper studies the impact of references on the initial responses generated by the LLM. There, we look at how sufficient correctness changes while we vary whether to give a reference to the LLM or not. Figure 11 illustrates the sufficient correctness of models from the Qwen3 family, as well as two frontier models: `gemini-2.5-pro` (Team et al., 2023) and `gpt-5.1` (gpt, 2025) on the FActScore Rare dataset. Although the frontier models perform better when a reference is given, we can see that Qwen3-4B is comparable to these frontier models when the reference text is given. Figure 12 and Figure 13 show the same comparison for the models from the Llama3 family and SmoLLM2 family, respectively. We observe that across model sizes and datasets, providing references consistently improves generation quality.

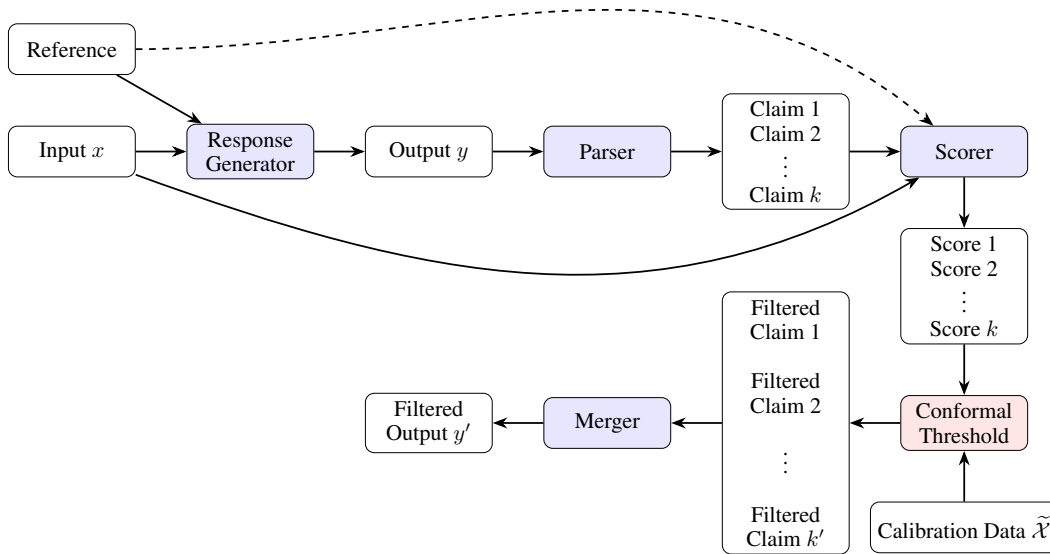


Figure 10: Given an input x and a reference text related to x , the Response Generator produces an output y , which is then parsed by the Parser into a list of claims. Each individual claim is subsequently scored by the Scorer, conditioned on the input x and, optionally, the reference text. These scores are passed to the conformal prediction algorithm, which filters out claims whose scores fall below a learned threshold. Finally, the remaining claims are merged into a single paragraph and returned to the user.

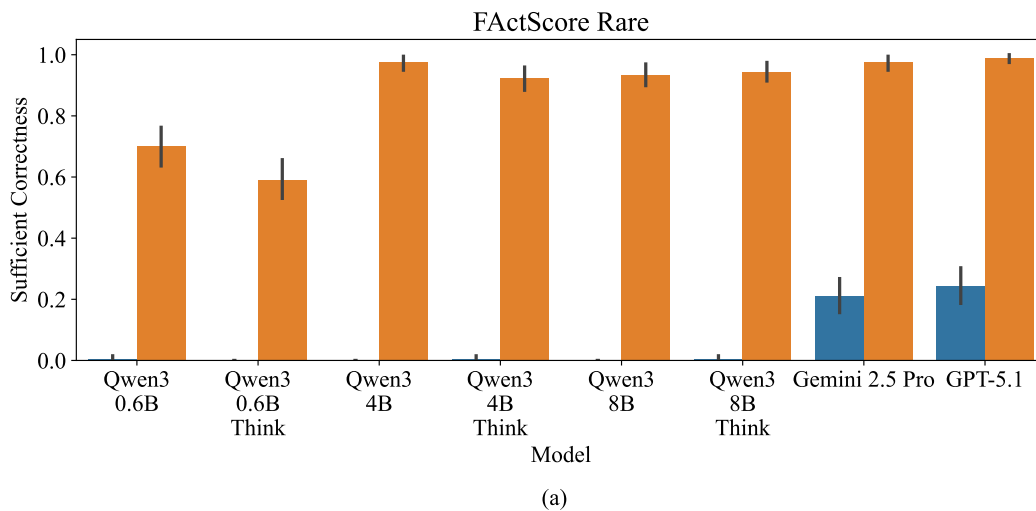


Figure 11: Sufficient correctness (SC) of Qwen3 and some frontier models on FActScore-Rare dataset, with and without access to references. Across these models, providing references consistently improves generation quality.

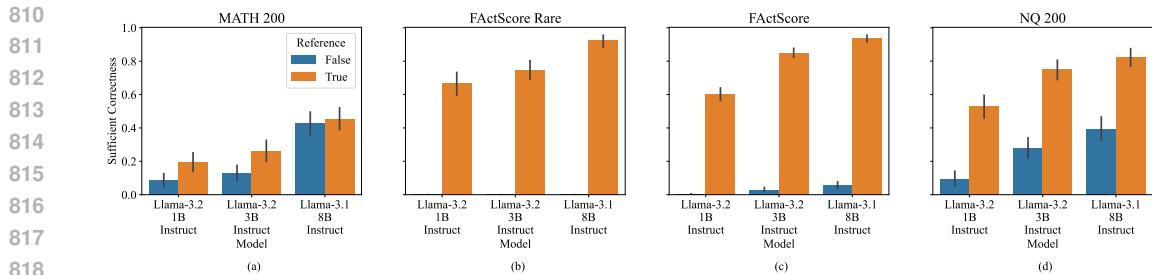


Figure 12: Sufficient correctness (SC) of Llama 3.1 and Llama 3.2 models on four datasets (MATH-200, FActScore-Rare, FActScore, NQ-200), with and without access to references. Across model sizes and datasets, providing references consistently improves generation quality.

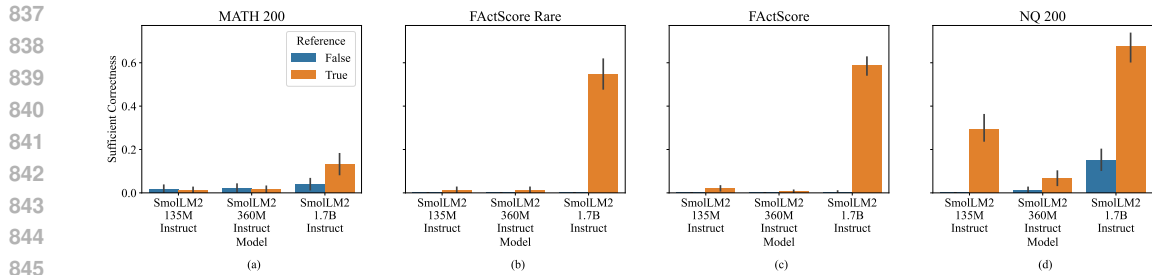


Figure 13: Sufficient correctness (SC) of SmoLLM2 models on four datasets (MATH-200, FActScore-Rare, FActScore, NQ-200), with and without access to references. Across model sizes and datasets, providing references consistently improves generation quality.

A.2.2 PROMPTING STRATEGIES FOR SCORERS

Figure 14 and Figure 15 show the overall performance comparison of model confidence score across different prompting strategies with different target factuality on the FActScore dataset and MATH-1K dataset, respectively.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

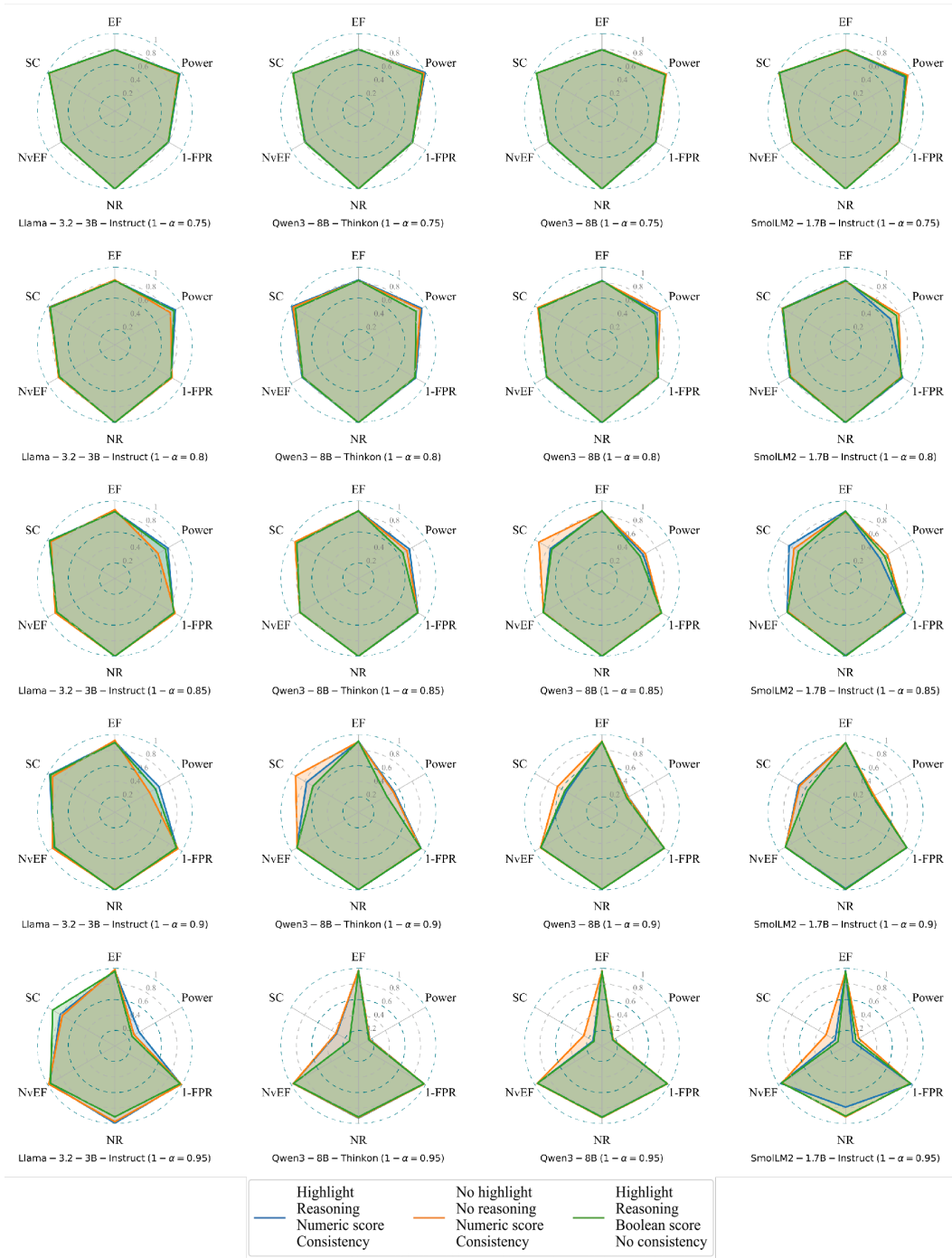


Figure 14: Overall performance comparison of model confidence scores across different prompting strategies with different target factuality ($1 - \alpha$) on FActScore dataset.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

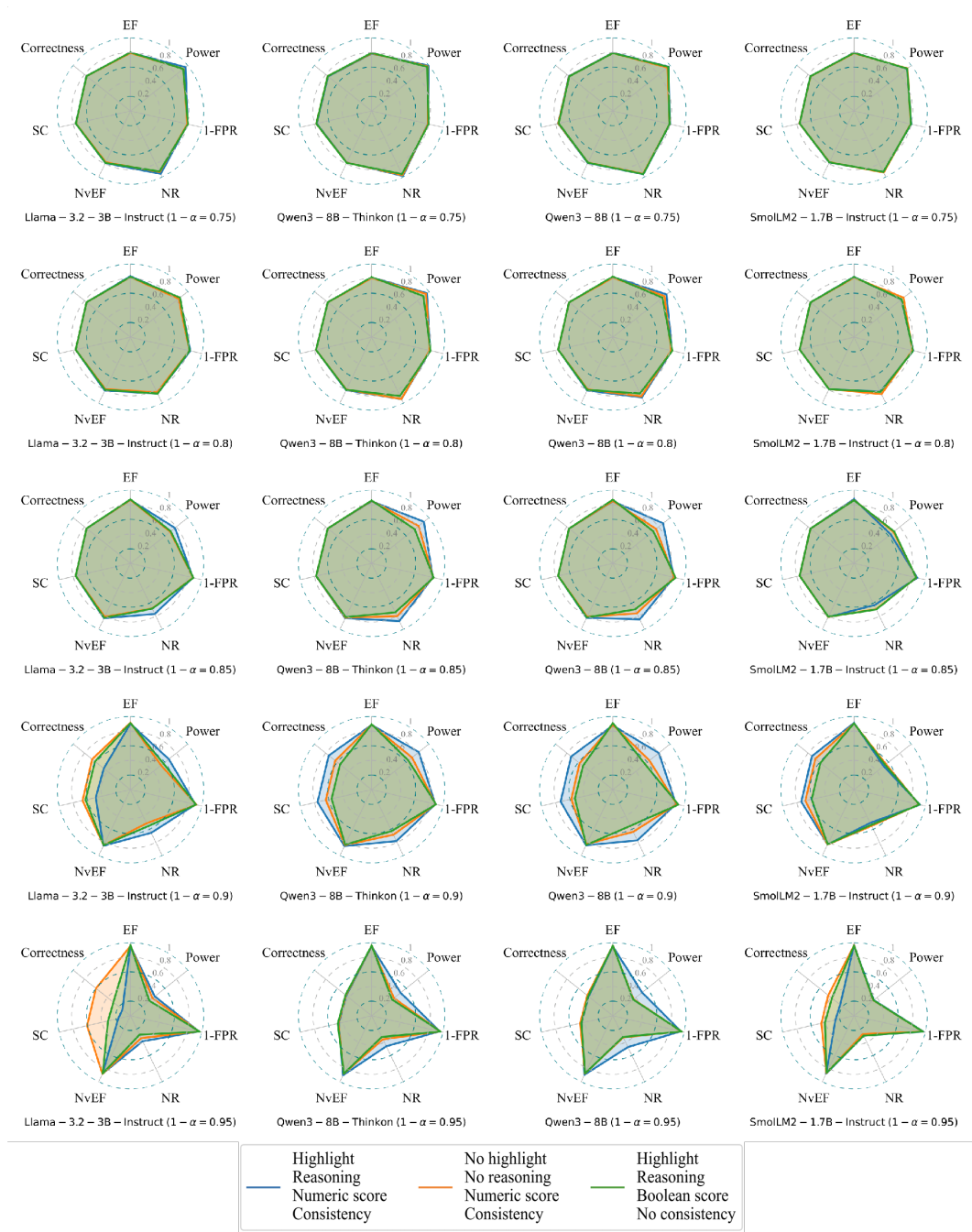
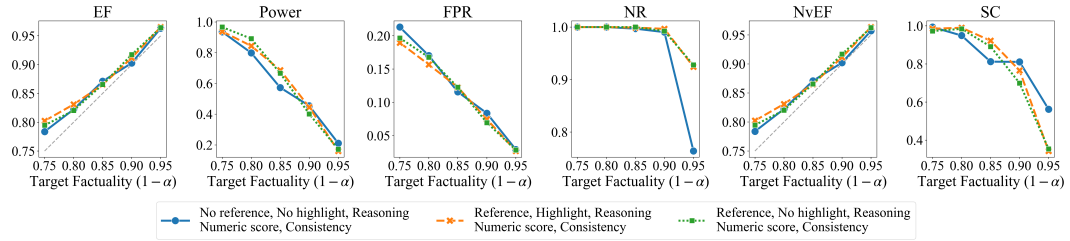


Figure 15: Overall performance comparison of model confidence scores across different prompting strategies with different target factuality ($1 - \alpha$) on MATH-1K dataset.

A.2.3 ROLE OF REFERENCES IN SCORING

To investigate the role of references to the model confidence score, we compare the case of giving and not giving the reference text to the scoring function under various metrics. Figure 16, Figure 4, and Figure 17 show the performance of various metrics using Qwen3-4B on the FActScore, MATH-1K, and NQ-1K datasets, respectively. Similarly, Figure 18, Figure 19, and Figure 20 show the performance of various metrics using Qwen3-8B on the FActScore, MATH-1K, and NQ-1K datasets.

972 Figure 21-23 and Figure 24-26 show the performances on the three datasets using Llama-3.2B
 973 and SmoLLM2-1.7B, respectively.
 974
 975
 976
 977



981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

Figure 16: Performance of model confidence score on FActScore with and without reference provided to scoring functions using gpt-5-nano and Qwen3-4B.

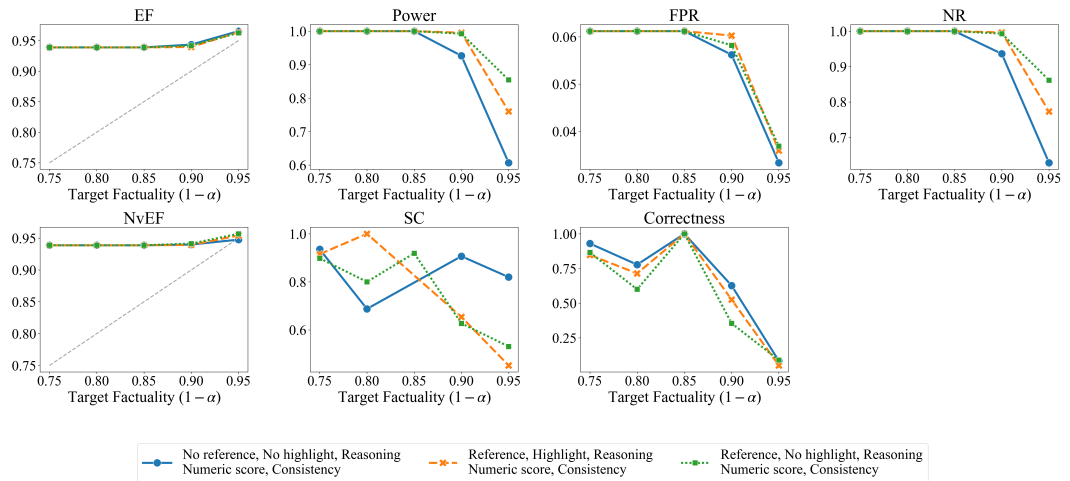


Figure 17: Performance of model confidence score on NQ-1K dataset with and without reference provided to scoring functions using gpt-5-nano and Qwen3-4B.

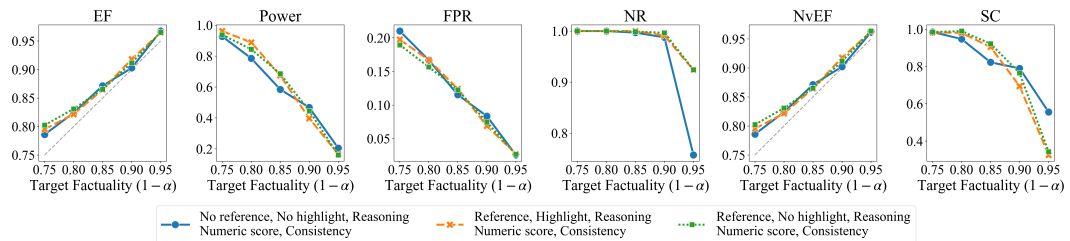


Figure 18: Performance of model confidence score on FActScore dataset with and without reference provided to scoring functions using gpt-5-nano and Qwen3-8B.

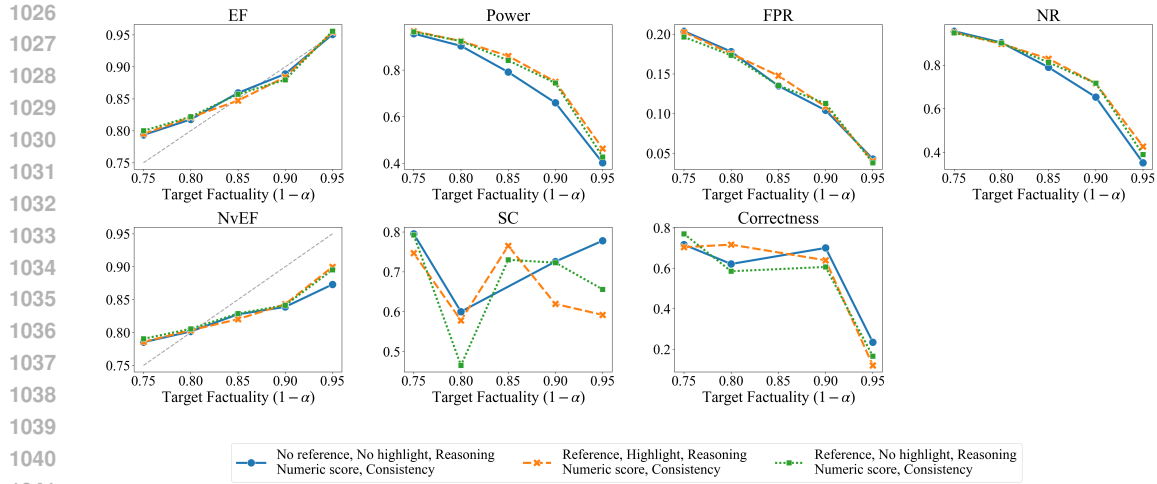


Figure 19: Performance of model confidence score on MATH-1K dataset with and without reference provided to scoring functions using gpt-5-nano and Qwen3-8B.

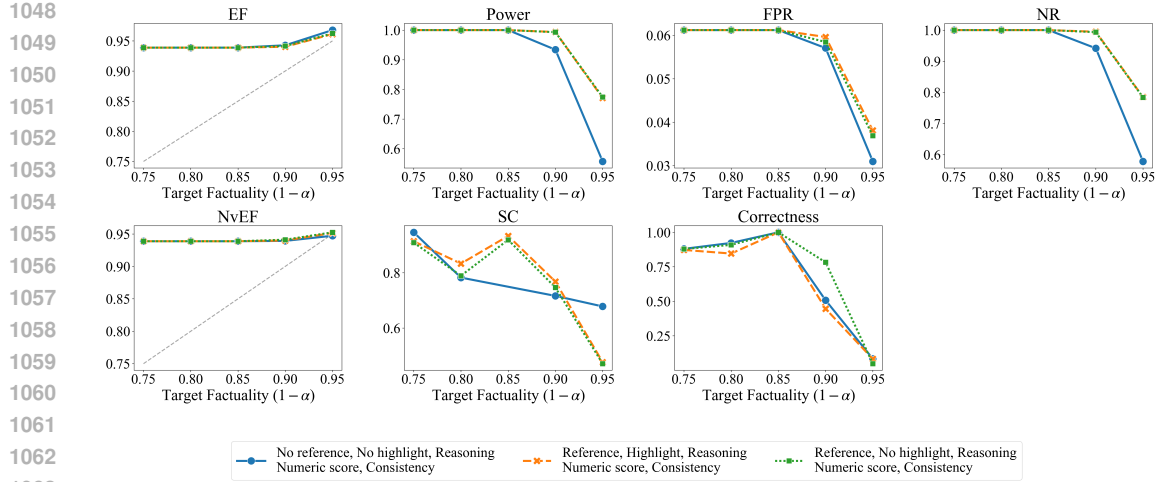


Figure 20: Performance of model confidence score on NQ-1K dataset with and without reference provided to scoring functions using gpt-5-nano and Qwen3-8B.

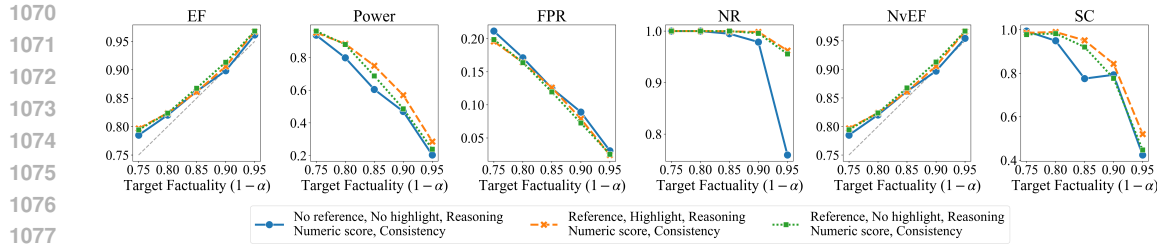


Figure 21: Performance of model confidence score on FActScore dataset with and without reference provided to scoring functions using gpt-5-nano and Llama-3.2-3B.

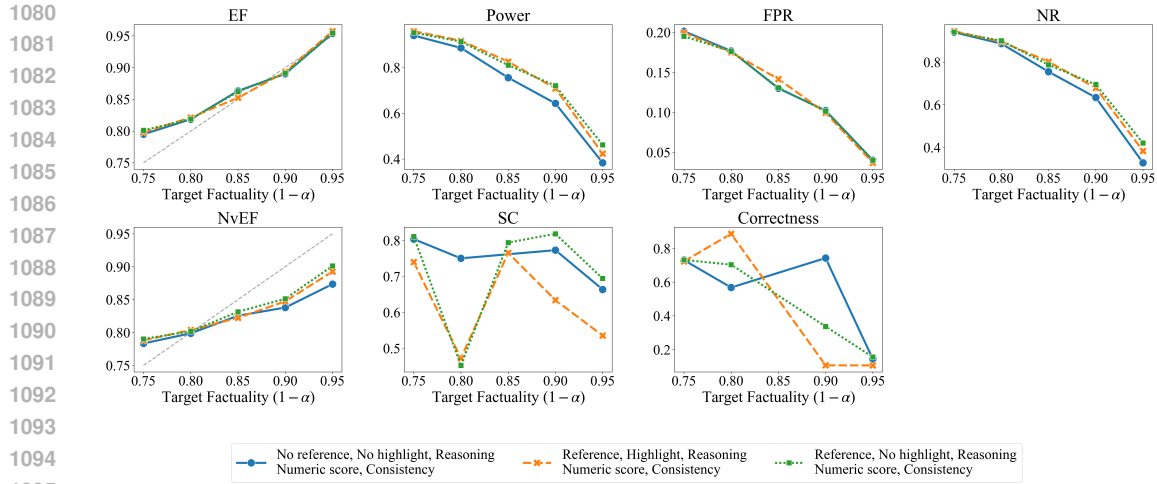


Figure 22: Performance of model confidence score on MATH-1K dataset with and without reference provided to scoring functions using gpt-5-nano and Llama-3.2-3B.

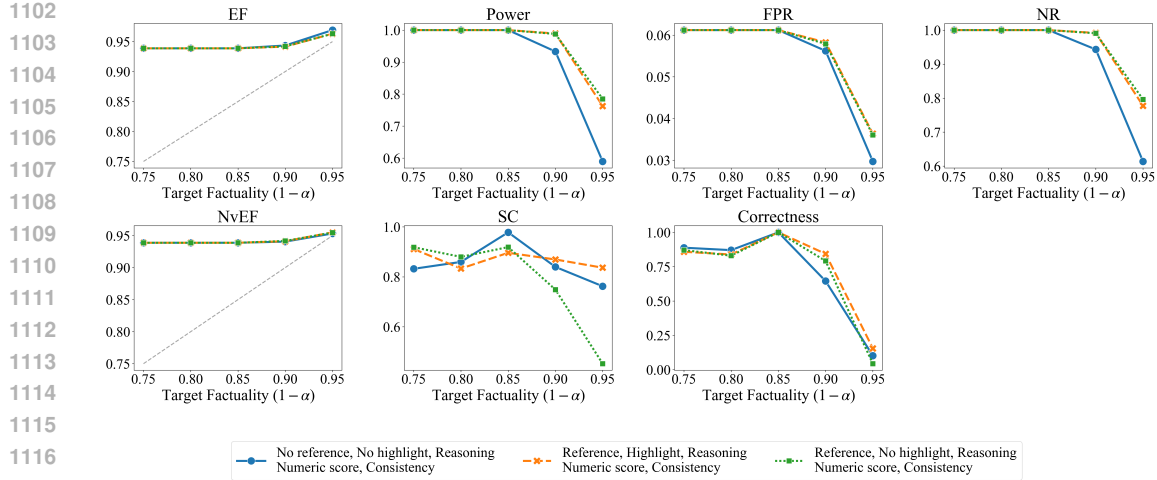


Figure 23: Performance of model confidence score on NQ-1K dataset with and without reference provided to scoring functions using gpt-5-nano and Llama-3.2-3B.

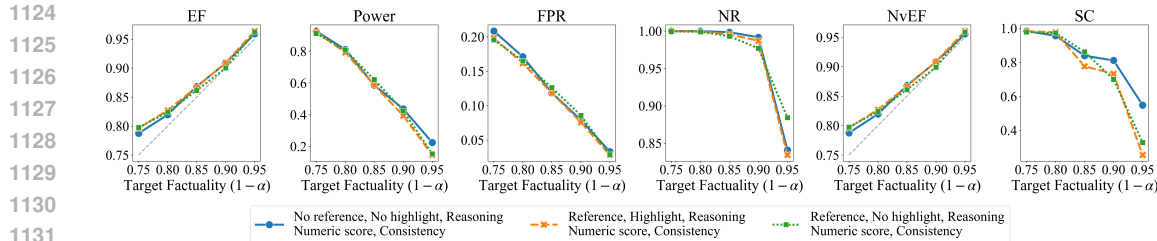


Figure 24: Performance of model confidence score on FActScore dataset with and without reference provided to scoring functions using gpt-5-nano and SmoLLM2-1.7B.

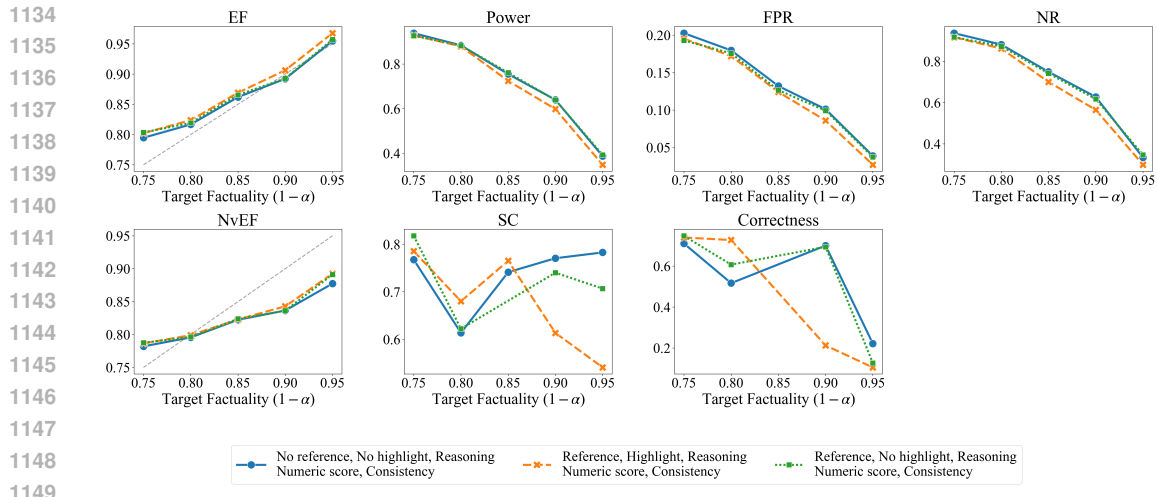


Figure 25: Performance of model confidence score on MATH-1K dataset with and without reference provided to scoring functions using gpt-5-nano and SmoLLM2-1.7B.

1152
1153
1154
1155
1156
1157
1158
1159

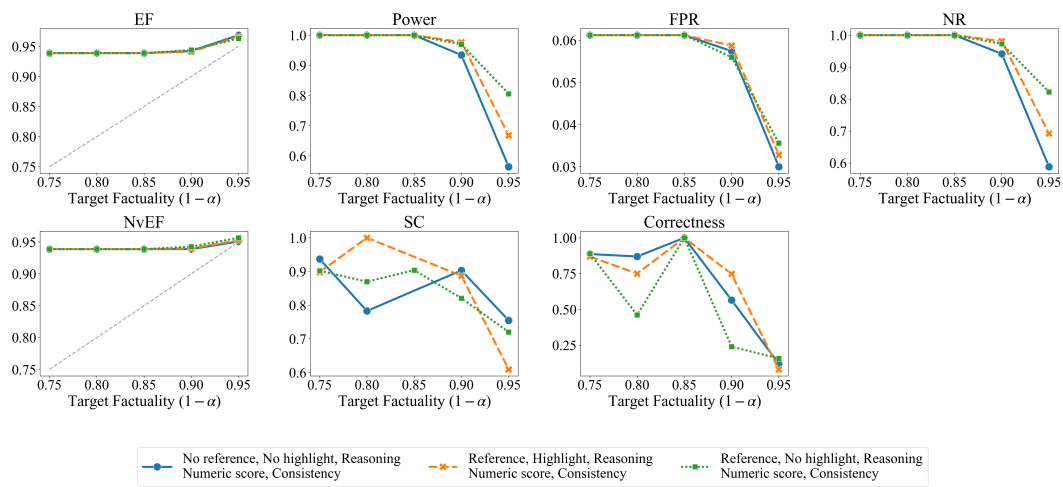


Figure 26: Performance of model confidence score on NQ-1K dataset with and without reference provided to scoring functions using gpt-5-nano and SmoLLM2-1.7B.

1178
1179
1180
1181

A.2.4 MODEL CHOICE FOR SCORERS

1184
1185
1186
1187

We compare different open-source models as scoring functions f , using gpt-5-nano as the generator G . All prompts include references, require highlighting and CoT, produce scalar scores, and apply consistency averaging. Figure 27 shows the assessment of how scorer model family and their scale affect factuality filtering with different target factuality $1 - \alpha$.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

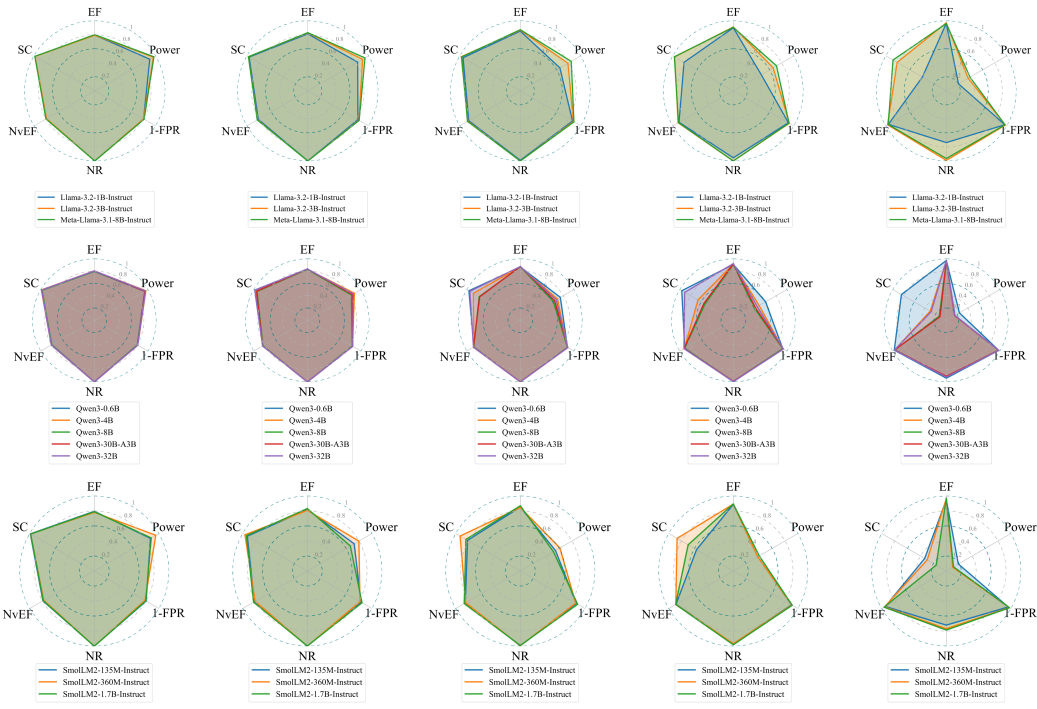


Figure 27: Overall performance comparison of model confidence scores across different model scales with different target factuality $1 - \alpha$.

A.2.5 ROBUSTNESS TO CALIBRATION DISTRIBUTION SHIFT

Although scoring functions differ in their design and accuracy, their real-world usefulness depends on robustness: can they maintain reliability under distribution shifts or adversarial perturbations? Figure 28 and 29 show that for the 3 datasets, when there is a distribution shift, the factuality guarantee no longer holds.

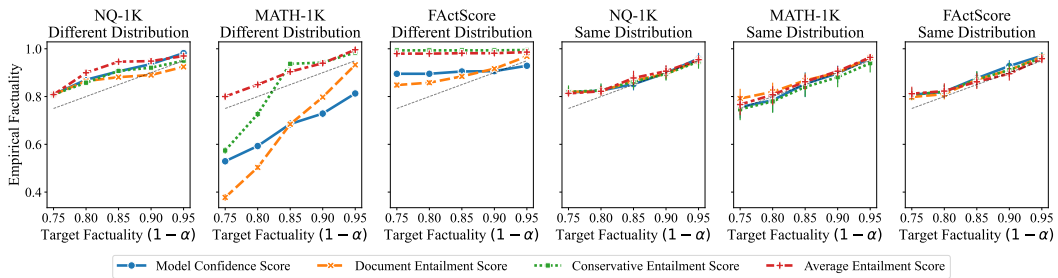


Figure 28: Empirical factuality (EF) on FactScore, MATH-1K, and NQ-1K, under two calibration settings: (i) calibration claims generated in Mohri & Hashimoto (2024), which comes from a different distribution (ii) calibration claims drawn from the same distribution as the test data. We use Llama3.2-3B as the scoring function. The results show how distribution shift in the calibration set affects conformal factuality guarantees.

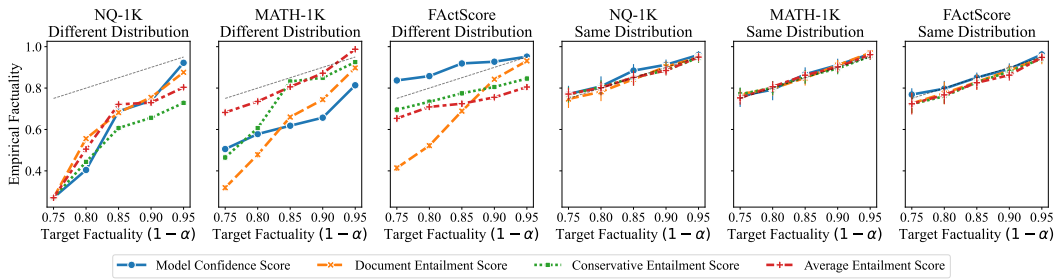


Figure 29: Empirical factuality (EF) on FActScore, MATH-1K, and NQ-1K, under two calibration settings: (i) calibration claims generated in Mohri & Hashimoto (2024), which comes from a different distribution (ii) calibration claims drawn from the same distribution as the test data. We use `SmolLM2-360M` as the scoring function. The results show how distribution shift in the calibration set affects conformal factuality guarantees.

A.2.6 ROBUSTNESS TO ADVERSARIAL DISTRACTORS

Figure 30 and Figure 31 show that for `Llama3.2-3B-Instruct` and `SmolLM2-1.7B-Instruct`, as the distractor rate increases, empirical factuality drops sharply on the FActScore dataset. Similar observations are also found on the MATH-1K and NQ-1K dataset, as shown in Figure 32-37 on MATH-1K and NQ-1K dataset.

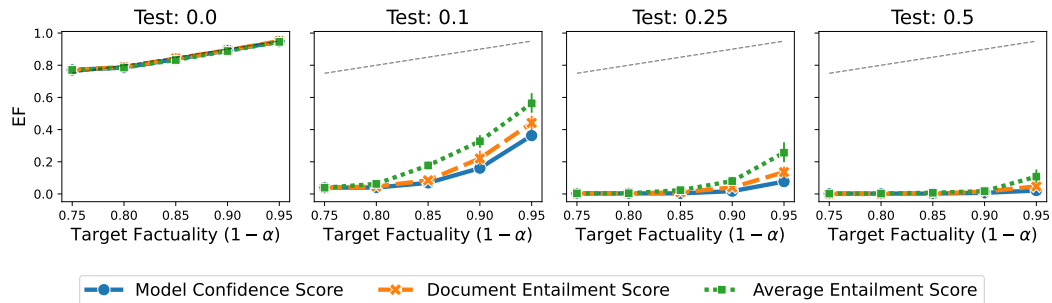


Figure 30: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with fixed clean calibration data. Three scoring methods compared on FActScore dataset with `Llama-3.2-3B-Instruct`.

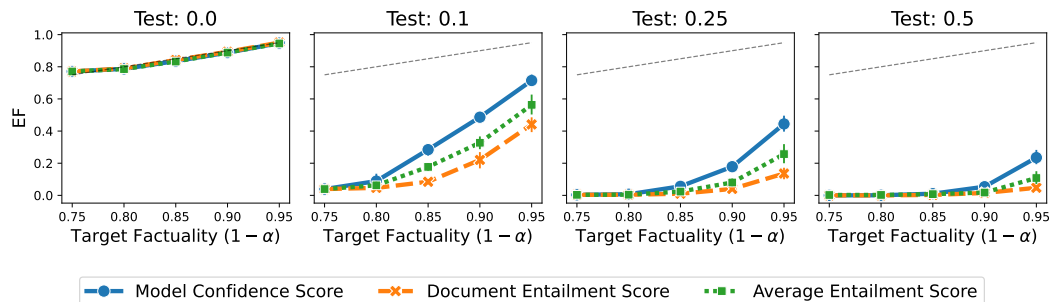
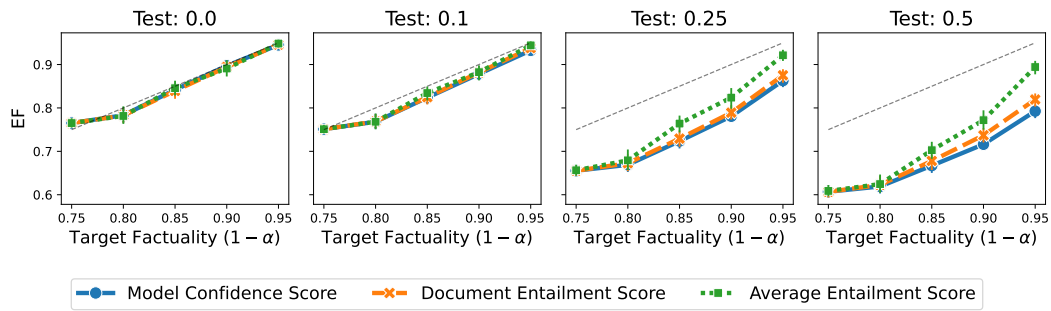
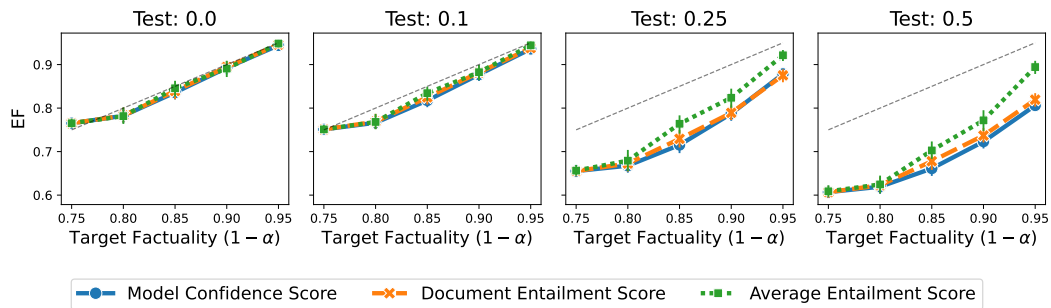


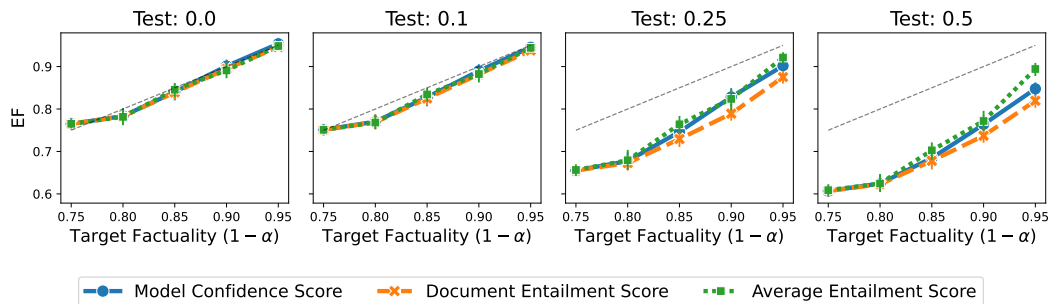
Figure 31: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with fixed clean calibration data. Three scoring methods compared on FActScore dataset with `SmolLM2-1.7B-Instruct`.



1307 Figure 32: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with fixed
 1308 clean calibration data. Three scoring methods compared on MATH-1K dataset with Qwen3-4B.



1327 Figure 33: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with
 1328 fixed clean calibration data. Three scoring methods compared on MATH-1K dataset with
 1329 Llama-3.2-3B-Instruct.



1348 Figure 34: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with
 1349 fixed clean calibration data. Three scoring methods compared on MATH-1K dataset with
 SmolLM2-1.7B-Instruct.

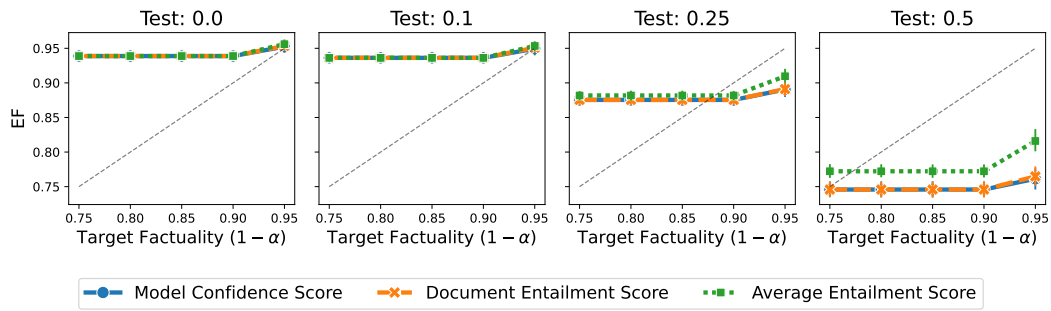


Figure 35: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with fixed clean calibration data. Three scoring methods compared on NQ-1K dataset with Qwen3-4B.

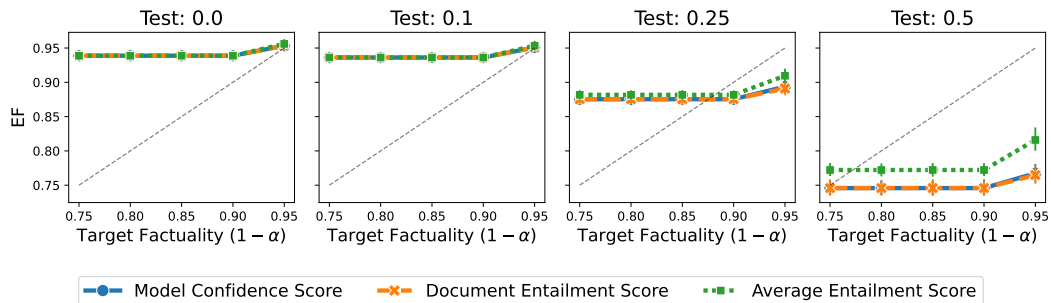


Figure 36: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with fixed clean calibration data. Three scoring methods compared on NQ-1K dataset with Llama-3.2-3B-Instruct.

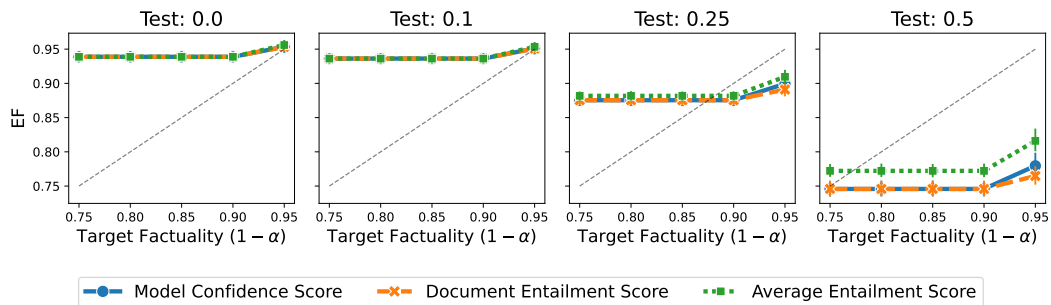


Figure 37: Empirical factuality under varying test hallucination proportions (0.0 to 0.5) with fixed clean calibration data. Three scoring methods compared on NQ-1K dataset with SmolLM2-1.7B-Instruct.

A.2.7 ROBUSTNESS WITH ADVERSARIAL PREPARED CALIBRATION

We show the result for a setting with a test set with a distractor proportion of 0.25 and varying the level of distractors in the calibration set. Figure 38-45 illustrates that for all three datasets, FActScore, MATH-1K, and NQ-1K, when the fractions of distractors are underestimated, the empirical factuality (EF) is still below the target EF. The EF reaches the target EF only when we introduce a large enough fraction of distractors to the calibration set. However, we note that this incurs a high cost on the non-empty rate.

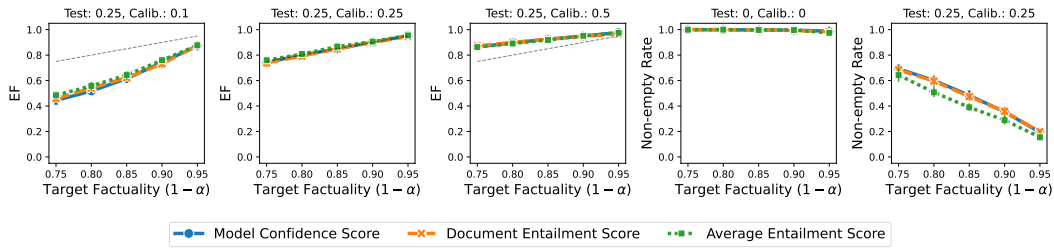


Figure 38: Empirical factuality and non-empty rates under different test-calibration distribution shifts. FActScore dataset with Llama-3.2-3B-Instruct. Shows trade-off between factuality guarantees and prediction coverage.

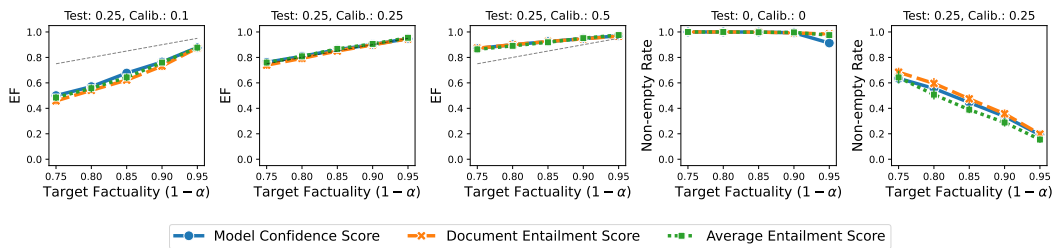


Figure 39: Empirical factuality and non-empty rates under different test-calibration distribution shifts. FActScore dataset with SmolLM2-1.7B-Instruct.

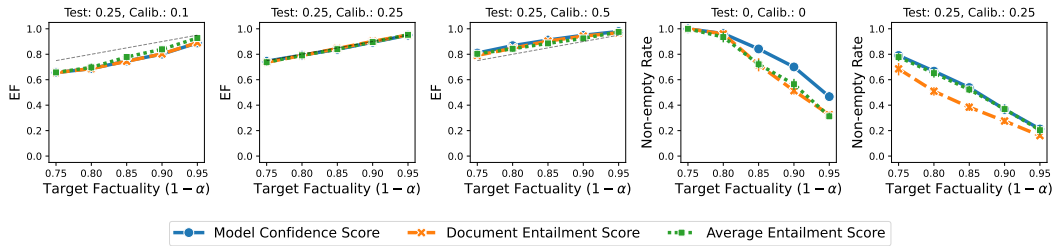


Figure 40: Empirical factuality and non-empty rates under different test-calibration distribution shifts. MATH-1K dataset with Qwen3-4B.

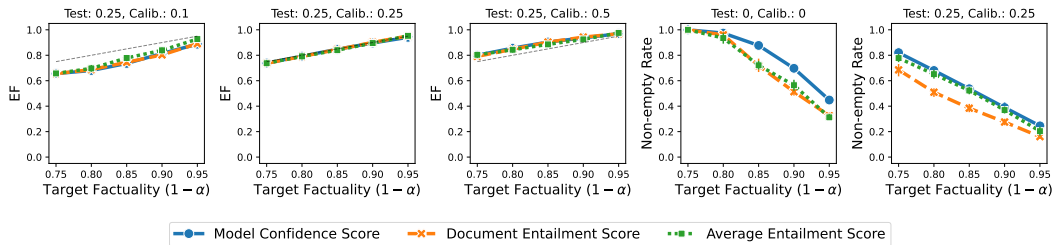


Figure 41: Empirical factuality and non-empty rates under different test-calibration distribution shifts. MATH-1K dataset with Llama-3.2-3B-Instruct.

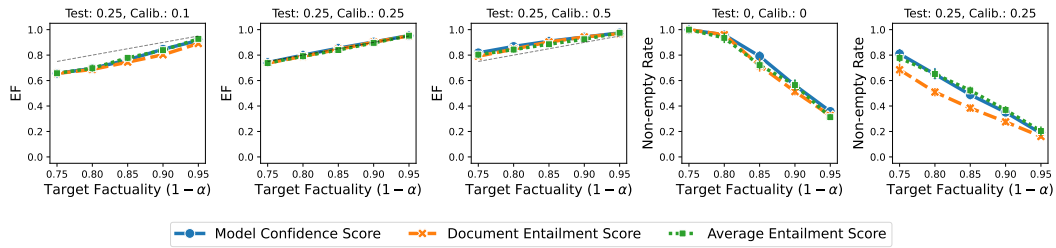


Figure 42: Empirical factuality and non-empty rates under different test-calibration distribution shifts. MATH-1K dataset with SmolLM2-1.7B-Instruct.

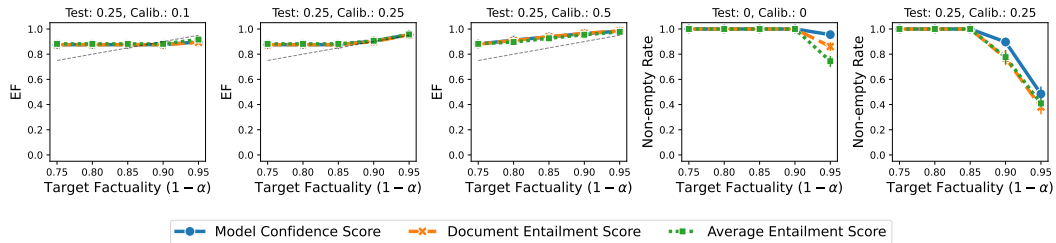


Figure 43: Empirical factuality and non-empty rates under different test-calibration distribution shifts. NQ-1K dataset with Qwen3-4B.

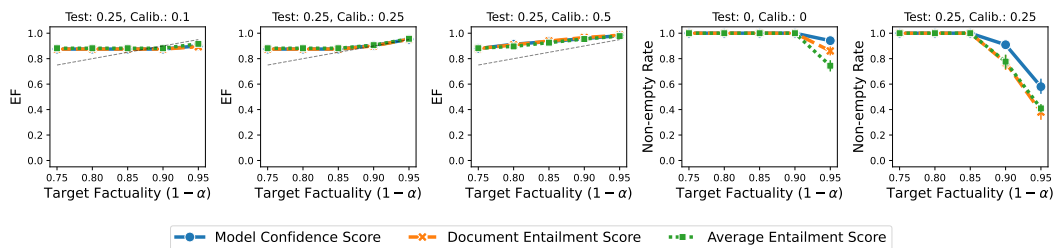


Figure 44: Empirical factuality and non-empty rates under different test-calibration distribution shifts. NQ-1K dataset with Llama-3.2-3B-Instruct.

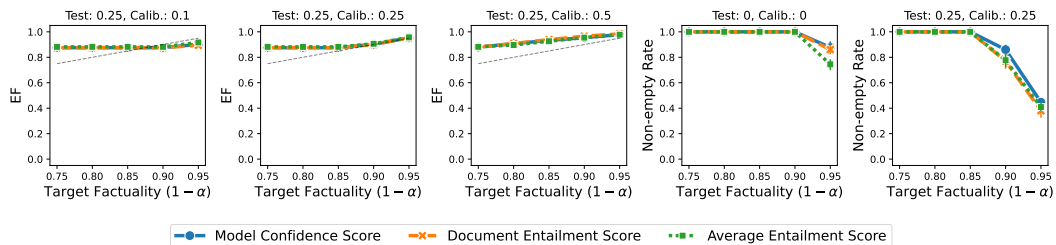


Figure 45: Empirical factuality and non-empty rates under different test-calibration distribution shifts. NQ-1K dataset with SmolLM2-1.7B-Instruct.

A.2.8 END-TO-END EVALUATION

While our earlier analyses isolate the contributions of individual components—such as references for generation, scoring function design, and robustness under distribution shifts—a complete understanding requires evaluating the entire pipeline as a whole. Component-level results alone can be

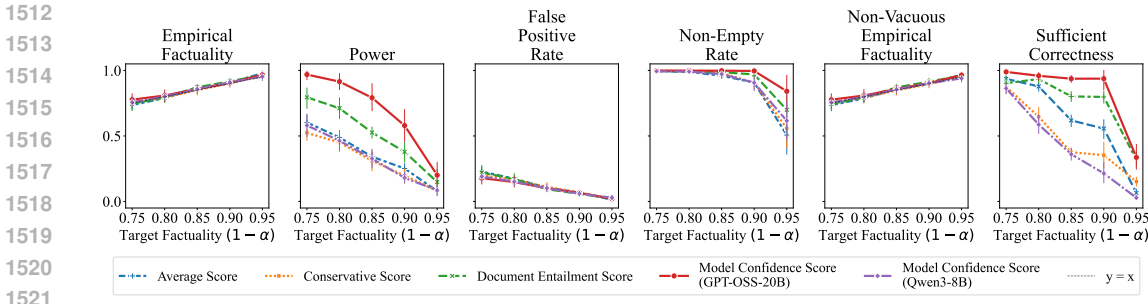


Figure 46: Performance of entailment- and confidence-based scorers using `gpt-oss-20b` and `Qwen3-8B` on FactScore.

misleading: a scorer may appear highly accurate in isolation but overly conservative when applied in practice, or references may improve raw generation quality but introduce new burdens on factuality filtering. Hence, now evaluate the full end-to-end pipeline. We use `gpt-oss-20b`, a mixture-of-experts model with 3.6B active parameters, as the response generator G and consider two options for the scoring function f : (i) using `gpt-oss-20b` as the scorer, which reflects the practical scenario in which the same model is available for generation and scoring; (ii) using `Qwen3-8B` as the scorer, which allows us to study whether a dense model can serve as an alternative to a FLOPs efficient MoE model for scoring.

We report results on FactScore, MATH, and NQ, evaluating both factuality and correctness. In addition, we analyze computational cost by comparing floating-point operation counts (FLOPs) required for inference across the two models.

Model	Total Tokens	Total FLOPs
<code>gpt-oss-20b</code> (3.6B active) (Agarwal et al., 2025)	2000	1.44×10^{13}
<code>Qwen3-8B</code> (8.19B active) (Yang et al., 2025)	2000	3.28×10^{13}
<code>DeepSeek-R1</code> (37B active) (DeepSeek-AI, 2025)	2000	1.5×10^{14}
<code>DeBERTa</code> (184M active) (He et al., 2020)	2000	4.9×10^{11}
<code>RoBERTa</code> (356M active) (Liu et al., 2019)	2000	1.6×10^{12}

Table 1: Estimated FLOPs for generating 1000 tokens with a 1000-token prompt (assuming KV caching).

Results and Discussion. Figure 46 shows that confidence-based scoring with `gpt-oss-20b` achieves higher power and sufficient correctness than with `Qwen3-8B`, especially at moderate target factuality levels. Despite activating fewer parameters, `gpt-oss-20b` consistently outperforms `Qwen3-8B`, suggesting advantages from the MoE architecture. Table 1 further highlights the efficiency gap: `gpt-oss-20b` requires less than half the FLOPs of `Qwen3-8B` for comparable workloads, yet delivers stronger factuality filtering. We also note that the document entailment score, which is based on `DeBERTa`, is almost more than two orders of magnitude efficient than `gpt-oss-20b`, performs similarly to the model confidence score on the non-empty rate and sufficient correctness.

Together, these results demonstrate that parameter-efficient MoE models can be both more effective and more economical than dense counterparts for end-to-end factuality evaluation.

A.3 DISTRACTORS

A major assumption of the conformal factuality framework is that the test datapoints are exchangeable with the calibration dataset. However, in a real-world scenario where we deploy the model, it is possible that this assumption is mildly violated.

To test how robust these scoring functions are under this scenario, for each query x_i , the corresponding reference text $R(x)$, and the set of claims $\{c_i\}$, we ask the LLM to modify each c_i , given x and

$R(x)$ such that it would become something that the model would hallucinate. We defer the prompts we used to generate these hallucination claims in Appendix A.5.6.

Our goal in creating these hallucination claims is that we want these claims to confuse the model so that the model would think that these claims are actually generated by them. To achieve this goal, after we generate a hallucinated claim, we ask the LLM to check if it thinks that the claim might be generated or hallucinated by itself (prompt in Appendix A.5.7), given the same $x, R(x)$. If the hallucination claim can cause the model to think that it is the one who generates it, given x and $R(x)$, then we keep this hallucination claim. Otherwise, we repeat this process and generate a new hallucination claim.

A.4 HUMAN EVALUATION

We conducted a human evaluation by randomly sampling 200 claims from the FActScore dataset. We had two students (referred to as A and B) label the factuality of these claims independently. We also used `gpt-5-nano` to label the same set of claims.

We then compared the agreement rates: `gpt-5-nano` vs. A, `gpt-5-nano` vs. B, and A vs. B. Table 2 shows the agreement rates. We found that the agreement rate between `gpt-5-nano` and a human labeler was around 75%, which was very similar to the inter-human agreement rate (A vs. B). This result suggests that using a capable LLM as a judge for factuality is a reasonable proxy, performing comparably to a human judge.

Pair	Labelers	Agreement Rate
Model-Human	<code>gpt-5-nano</code> vs. Student A	76.5%
Model-Human	<code>gpt-5-nano</code> vs. Student B	77.0%
Human-Human	Student A vs. Student B	73.0%

Table 2: Agreement rates on factuality labels between `gpt-5-nano` and two human annotators, as well as between the two human annotators themselves.

A.5 PROMPTS

A.5.1 GENERATOR (WITH REFERENCE)

```
You are a helpful assistant that answers queries strictly based on the
provided reference text.

Instructions:
1. You will be given:
  - A reference text
  - A query
2. Use only the information from the reference text to answer the query.
3. Do not include any information not supported by the reference text.

Output Requirements:
- Output ONLY a single VALID JSON5 object with EXACTLY this schema:
  {
    "response": "...answer strictly based on the reference text..."
  }

JSON5 Rules:
- Use DOUBLE QUOTES (") for all keys and all string values.
- Escape double quotes inside string values as \".
- Escape backslashes as \\.
- No trailing commas in objects or arrays.
- Use the exact top-level container specified and close it properly.
- Do not include comments, code fences, or any text outside the JSON5
  output.
- Follow the schema exactly; do not add or omit keys.
```

1620 Do NOT include:
1621 - Any text, explanations, comments, or formatting outside of the JSON5.
1622 - Any code block delimiters (e.g., ``json).
1623
1624 Input:
1625 Reference Text: {reference}
1626 Query: {query}
1627
1628 (Reiteration of the instruction)
1629 Answer the query strictly using only the reference text, and return a
1630 single JSON5 object with the key "response" only.
1631
1632 Output:

1633 A.5.2 GENERATOR (WITHOUT REFERENCE)

1635 You are a helpful assistant that answers queries.
1636
1637 Instructions:
1638 1. You will be given:
1639 - A query
1640
1641 Output Requirements:
1642 - Output ONLY a single VALID JSON5 object with EXACTLY this schema:
1643 {
1644 "response": "...answer..."
1645 }
1646
1647 JSON5 Rules:
1648 - Use DOUBLE QUOTES (") for all keys and all string values.
1649 - Escape double quotes inside string values as \".
1650 - Escape backslashes as \\
1651 - No trailing commas in objects or arrays.
1652 - Use the exact top-level container specified and close it properly.
1653 - Do not include comments, code fences, or any text outside the JSON5
1654 output.
1655 - Follow the schema exactly; do not add or omit keys.
1656
1657 Do NOT include:
1658 - Any text, explanations, comments, or formatting outside of the JSON5.
1659 - Any code block delimiters (e.g., ``json).
1660
1661 Input:
1662 Query: {query}
1663
1664 (Reiteration of the instruction)
1665 Answer the query and return a single JSON5 object with the key "response"
1666 only.
1667
1668 Output:

1665 A.5.3 PARSER

1667 You are an AI assistant tasked with breaking down input text into small,
1668 self-contained claims for easy human verification.
1669
1670 Instructions:
1671 1. Parse the provided text into concise, independent, and non-overlapping
1672 subclaims.
1673 2. Ensure each subclaim is:
- As small and specific as possible.
- Independent and self-contained.

1674 - Do not use pronouns like he, she, his, her, it, its, etc.
1675 - Explicitly mention subjects.
1676 - Factually complete without relying on context from other subclaims.
1677 3. If the provided text is not a full sentence, use the provided text
1678 verbatim as the subclaim

1679 Output Requirements:
1680 1. The result must be a VALID and COMPLETE JSON list of dictionaries.
1681 2. Each dictionary must have the following structure:
1682 {
1683 "subclaim": "Subclaim text"
1684 }
1685

1685 JSON Rules:
1686 - Ensure the JSON is STRICTLY VALID:
1687 - Use DOUBLE QUOTES (") for all keys and string values.
1688 - DO NOT include trailing commas after the LAST item in arrays or
1689 objects.
1690 - Ensure ALL dictionaries are enclosed in curly braces {}.
1691 - Ensure the JSON list is ENCLOSED in square brackets [].
1692 - CLOSE the JSON list properly with a closing square bracket].
1693 - DO NOT include any code block delimiters.
1694 - DO NOT include any additional text, explanations, or comments in the
1695 outputreturn ONLY the JSON list.
1696

1696 Critical Requirement:
1697 - Validate and ensure the JSON is complete, properly CLOSED, and
1698 parseable. The final JSON MUST pass validation without errors.
1699

1698 (Reiteration of the instruction)
1699 1. Parse the provided text into concise, independent, and non-overlapping
1700 subclaims.
1701 2. Ensure each subclaim is:
1702 - As small and specific as possible.
1703 - Independent and self-contained.
1704 - Do not use pronouns like he, she, his, her, it, its, etc.
1705 - Explicitly mention subjects.
1706 - Factually complete without relying on context from other subclaims.
1707 3. If the provided text is not a full sentence, use the provided text
1708 verbatim as the subclaim

1708 Input:
1709 {input}
1710

1710 Output:
1711

1712 A.5.4 LABELER (WITHOUT GROUND-TRUTH ANSWER)

1714 You are an AI assistant tasked with assigning a label to a claim based
1715 on its factuality.
1716

1717 Instructions:
1718 1. You are given:
1719 - A reference text for the query.
1720 - A query.
1721 - A claim made in response to the query.
1722 2. Rate the factuality of the claim as a boolean:
1723 - true -> The claim is well-supported by the reference text.
1724 - false -> The claim contradicts the reference text, or is not well-
1725 supported by it.
1726 - If it is ambiguous, answer False.
1727 3. Provide:
1728 - The part(s) of the reference text that directly support your
1729 decision.
1730 - A reasoning statement describing your rationale.

1728 4. You must assign either true or false. Never return null or None.
1729

1730 Output Requirements:
1731 - Output ONLY a single VALID JSON5 object with EXACTLY these keys:
1732 {
1733 "highlighted_text": "Part(s) of the reference text that support the
1734 decision.",
1735 "reasoning": "A reasoning statement describing your rationale.",
1736 "answer": true
1737 }
1738 - "answer" must be a boolean (true/false).

1739 JSON5 Rules:
1740 - Use DOUBLE QUOTES (") for all keys and all string values.
1741 - Escape double quotes inside string values as \".
1742 - Escape backslashes as \\
1743 - No trailing commas in objects or arrays.
1744 - Use the exact top-level container specified and close it properly.
1745 - Do not include comments, code fences, or any text outside the JSON5
1746 output.
1747 - Follow the schema exactly; do not add or omit keys.

1748 Do NOT include:
1749 - Any text, explanations, comments, or formatting outside of the JSON5.
1750 - Any code block delimiters (e.g., ``json).

1751 Examples:
1752 Example Input:
1753 Reference Text: Michael Scott is a fictional character in the NBC sitcom
1754 The Office, portrayed by Steve Carell. Michael is the regional
1755 manager of the Scranton, Pennsylvania branch of Dunder Mifflin, a
1756 paper company, for the majority of the series. Like his counterpart
1757 in the original British version of the show, David Brent, he is
1758 characterized as a largely incompetent, unproductive, unprofessional
1759 boss, though he is depicted as kinder and occasionally shown to be
1760 effective at his job in key moments.
1761 Query: Tell me a paragraph bio of Michael Scott.
1762 Claim: The fictional character Michael Scott is the regional manager of a
1763 paper company.

1764 Example Output:
1765 {
1766 "highlighted_text": "Michael Scott is a fictional character in the NBC
1767 sitcom The Office, portrayed by Steve Carell. Michael is the regional
1768 manager of the Scranton, Pennsylvania branch of Dunder Mifflin, a
1769 paper company, for the majority of the series.",
1770 "reasoning": "Reference explicitly states Michael is regional manager
1771 at a paper company.",
1772 "answer": true
1773 }

1774 Example Input:
1775 Reference Text: Michael Scott is a fictional character in the NBC sitcom
1776 The Office, portrayed by Steve Carell. Michael is the regional
1777 manager of the Scranton, Pennsylvania branch of Dunder Mifflin, a
1778 paper company, for the majority of the series. Like his counterpart
1779 in the original British version of the show, David Brent, he is
1780 characterized as a largely incompetent, unproductive, unprofessional
1781 boss, though he is depicted as kinder and occasionally shown to be
1782 effective at his job in key moments.
1783 Query: Tell me a paragraph bio of Michael Scott.
1784 Claim: The portrayal of Michael Scott in the NBC sitcom The Office is
1785 similar to that of David Brent, in the British version.
1786 Example Output:

```

1782 {
1783   "highlighted_text": "Like his counterpart in the original British
1784     version of the show, David Brent, he is characterized as a largely
1785     incompetent, unproductive, unprofessional boss, though he is depicted
1786     as kinder and occasionally shown to be effective at his job in key
1787     moments.",
1788   "reasoning": "Reference compares Michael Scotts characterization to
1789     David Brents, indicating similarity.",
1790   "answer": true
1791 }
1792 Example Input:
1793 Reference Text: Michael Scott is a fictional character in the NBC sitcom
1794   The Office, portrayed by Steve Carell. Michael is the regional
1795   manager of the Scranton, Pennsylvania branch of Dunder Mifflin, a
1796   paper company, for the majority of the series. Like his counterpart
1797   in the original British version of the show, David Brent, he is
1798   characterized as a largely incompetent, unproductive, unprofessional
1799   boss, though he is depicted as kinder and occasionally shown to be
1800   effective at his job in key moments.
1801 Query: Tell me a paragraph bio of Michael Scott.
1802 Claim: Michael Scott is the founder of The Michael Scott Paper Company.
1803 Example Output:
1804 {
1805   "highlighted_text": "Michael is the regional manager of the Scranton,
1806     Pennsylvania branch of Dunder Mifflin, a paper company, for the
1807     majority of the series.",
1808   "reasoning": "Reference mentions only Dunder Mifflin; founding another
1809     company is not supported.",
1810   "answer": false
1811 }
1812 Example Input:
1813 Reference Text: Michael Scott is a fictional character in the NBC sitcom
1814   The Office, portrayed by Steve Carell. Michael is the regional
1815   manager of the Scranton, Pennsylvania branch of Dunder Mifflin, a
1816   paper company, for the majority of the series. Like his counterpart
1817   in the original British version of the show, David Brent, he is
1818   characterized as a largely incompetent, unproductive, unprofessional
1819   boss, though he is depicted as kinder and occasionally shown to be
1820   effective at his job in key moments.
1821 Query: Tell me a paragraph bio of Michael Scott.
1822 Claim: Michael Scott is the CEO of Dunder Mifflin.
1823 Example Output:
1824 {
1825   "highlighted_text": "Michael is the regional manager of the Scranton,
1826     Pennsylvania branch of Dunder Mifflin, a paper company, for the
1827     majority of the series.",
1828   "reasoning": "Reference states regional manager, not CEO.",
1829   "answer": false
1830 }
1831 Input:
1832 Reference Text: {reference}
1833 Query: {query}
1834 Claim: {claim}
1835 (Reiteration of the instruction)
1836 Return a single JSON5 object with "highlighted_text", "reasoning", and "
1837   answer" (true if supported; false if contradicted or unsupported).
1838   Assign true or falsenever null/None.
1839 Output:

```

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

A.5.5 LABELER (WITH GROUND-TRUTH ANSWER)

You are an AI assistant tasked with assigning a label to a claim based on its factuality.

Instructions:

1. You are given:
 - A reference text for the query.
 - A query.
 - A provided solution (final answer)
 - A claim made in response to the query.
2. Rate the factuality of the claim as a boolean:
 - true The claim is well-supported by the reference text or match the given provided solution (final answer).
 - false The claim contradicts the reference text or is not well-supported by it and contradicts the provided solution (final answer).
3. Provide:
 - The part(s) of the reference text or the provided solution that directly support your decision.
 - A reasoning statement describing your rationale.
4. You must assign either true or false. Never return null or None.

Output Requirements:

- Output ONLY a single VALID JSON5 object with EXACTLY this schema:


```
{
  "highlighted_text": "Part(s) of the reference text or the provided
  solution that directly support the decision.",
  "reasoning": "A reasoning statement describing your rationale.",
  "answer": true
}
```
- "answer" must be a boolean (true/false).

JSON5 Rules:

- Use DOUBLE QUOTES (") for all keys and all string values.
- Escape double quotes inside string values as \".
- Escape backslashes as \\\.
- No trailing commas in objects or arrays.
- Use the exact top-level container specified and close it properly.
- Do not include comments, code fences, or any text outside the JSON5 output.
- Follow the schema exactly; do not add or omit keys.

Do NOT include:

- Any text, explanations, comments, or formatting outside of the JSON5.
- Any code block delimiters (e.g., ```json).

Examples:

Example Input:

Reference Text: "Paris is the capital of France."

Answer: "Paris"

Query: "What is the capital of France?"

Claim: "Paris is the capital of France."

Example Output:

```
{
  "reasoning": "Claim matches reference and provided solution.",
  "answer": true
}
```

Example Input:

Reference Text: "Michael is the regional manager of Dunder Mifflin."

Answer: "Regional manager"

```

1890 Query: "What is Michaels role at Dunder Mifflin?"
1891 Claim: "Michael is the CEO of Dunder Mifflin."
1892
1893 Example Output:
1894 {
1895   "reasoning": "Claim contradicts the reference and solution.",
1896   "answer": false
1897 }
1898
1898 Example Input:
1899 Reference Text: "Compute 12 3 + 4."
1900 Answer: "40"
1901 Query: "Evaluate 12 3 + 4."
1902 Claim: "The value is 40."
1903
1903 Example Output:
1904 {
1905   "reasoning": "Computation equals 40; claim matches the solution.",
1906   "answer": true
1907 }
1908
1908 Example Input:
1909 Reference Text: "Compute 12 3 + 4."
1910 Answer: "40"
1911 Query: "Evaluate 12 3 + 4."
1912 Claim: "The value is 42."
1913
1913 Example Output:
1914 {
1915   "reasoning": "The solution gives 40, not 42.",
1916   "answer": false
1917 }
1918
1918 Example Input:
1919 Reference Text: "The sum of the interior angles of a triangle is 180
1920 degrees."
1921 Answer: "180 degrees"
1922 Query: "What is the sum of the interior angles of a triangle?"
1923 Claim: "The sum is 180 degrees."
1924
1924 Example Output:
1925 {
1926   "reasoning": "Claim agrees with the stated geometric fact and the
1927   solution.",
1928   "answer": true
1929 }
1930
1930 Input:
1931 Reference Text: {reference}
1932 Answer: {answer}
1933 Query: {query}
1934 Claim: {claim}
1935
1935 (Reiteration of the instruction)
1936 Return a single JSON5 object with "highlighted_text", "reasoning", and "
1937 answer" (true if supported; false if contradicted or unsupported).
1938 Assign true or falsenever null/None.
1939
1939 Output:

```

1942 A.5.6 ATTACKER
1943

1944 You are a helpful assistant that modifies a correct claim into a false
1945 but plausible claim (hallucination) given a reference text and a
1946 query.
1947

1948 Instructions:
1949 1. Read the reference text, the query, the original claim, and the two
1950 lists of prior modified claims (that do not work / that do work).
1951 2. Generate exactly one false claim that:
1952 - Directly answers the query.
1953 - Contradicts or is not supported by the reference text.
1954 - Appears plausible given the reference and query (e.g., near-miss
1955 number/date, closely related entity, role swap).
1956 - Does not match or paraphrase any claim in either provided list.
1957 3. Provide a reasoning statement describing why the false claim might
1958 look plausible.

1959 Output Requirements:
1960 - Output ONLY a single VALID JSON5 object with EXACTLY this schema:
1961 {
1962 "reasoning": "A reasoning statement describing why the false claim
1963 might look plausible.",
1964 "subclaim": "One-sentence false but plausible answer to the query."
1965 }
1966 - The "subclaim" must be a standalone sentence that answers the query and
1967 is false with respect to the reference.

1968 JSON5 Rules:
1969 - Use DOUBLE QUOTES (") for all keys and all string values.
1970 - Escape double quotes inside string values as \".
1971 - Escape backslashes as \\
1972 - No trailing commas in objects or arrays.
1973 - Use the exact top-level container specified and close it properly.
1974 - Do not include comments, code fences, or any text outside the JSON5
1975 output.
1976 - Follow the schema exactly; do not add or omit keys.

1977 Do NOT include:
1978 - Any text, explanations, comments, or formatting outside of the JSON5.
1979 - Any code block delimiters (e.g., ``json``).

1980 Examples:
1981 Example Input:
1982 Modified claims that do not work: []
1983 Modified claims that do work: []
1984 Reference Text: "Gustave Eiffel oversaw the construction of the Eiffel
1985 Tower, which opened in 1889. The towers design is credited to Gustave
1986 Eiffel."
1987 Query: "Who designed the Eiffel Tower?"
1988 Claim: "Gustave Eiffel designed the Eiffel Tower."

1989 Example Output:
1990 {
1991 "reasoning": "Reference credits Eiffel; another associated architect
1992 name may look credible.",
1993 "subclaim": "Stephen Sauvestre designed the Eiffel Tower."
1994 }
1995

1996 Example Input:
1997 Modified claims that do not work: []
1998 Modified claims that do work: []
1999 Reference Text: "The capital of Australia is Canberra."
2000 Query: "What is the capital of Australia?"
2001 Claim: "Canberra."

2002 Example Output:

```

1998 {
1999   "reasoning": "Larger city misconception is common.",
2000   "subclaim": "Sydney is the capital of Australia."
2001 }
2002
2003 Input:
2004 Modified claims that do not work: {claims_not_work}
2005 Modified claims that do work: {claims_work}
2006 Reference Text: {reference}
2007 Query: {query}
2008 Claim: {claim}
2009
2010 (Reiteration of the instruction)
2011 Generate exactly one false but plausible claim answering the query,
2012 unsupported by or contradicting the reference, not overlapping with
2013 provided lists; then return a JSON5 object with "reasoning" and "
2014 subclaim".
2015
2016 Output:

```

A.5.7 CONFUSEE

```

2017 You are a helpful assistant that evaluates whether a claim sounds like
2018 something you (an LLM) would plausibly say and potentially
2019 hallucinate based on the provided reference text and query.
2020
2021 Instructions:
2022 1. Read the reference text, the query, and the claim.
2023 2. Decide whether the claim is a plausible LLM-style statement that could
2024 arise from the reference and query even if it is not supported or is
2025 false.
2026 3. Return true if the claim is plausibly hallucinated; otherwise return
2027 false.
2028 4. Provide a reasoning statement describing your rationale.
2029
2030 Output Requirements:
2031 - Output ONLY a single VALID JSON5 object with EXACTLY this schema:
2032 {
2033   "reasoning": "A reasoning statement describing your rationale.",
2034   "answer": (true|false)
2035 }
2036 - "answer" must be a boolean (true/false).
2037
2038 JSON5 Rules:
2039 - Use DOUBLE QUOTES (") for all keys and all string values.
2040 - Escape double quotes inside string values as \".
2041 - Escape backslashes as \\.
2042 - No trailing commas in objects or arrays.
2043 - Use the exact top-level container specified and close it properly.
2044 - Do not include comments, code fences, or any text outside the JSON5
2045 output.
2046 - Follow the schema exactly; do not add or omit keys.
2047
2048 Do NOT include:
2049 - Any text, explanations, comments, or formatting outside of the JSON5.
2050 - Any code block delimiters (e.g., ``json).
2051
2052 Examples:
2053 Example Input:
2054 Reference Text: "The capital of Australia is Canberra."
2055 Query: "What is the capital of Australia?"
2056 Claim: "Sydney is the capital of Australia."
2057
2058 Example Output:

```

```

2052 {
2053   "reasoning": "Common misconception favors the largest city over the
2054   capital.",
2055   "answer": true
2056 }
2057 Example Input:
2058 Reference Text: "Paris is the capital of France."
2059 Query: "What is the capital of France?"
2060 Claim: "Paris is the capital of France."
2061 Example Output:
2062 {
2063   "reasoning": "Claim is directly supported by the reference; not a
2064   hallucination.",
2065   "answer": false
2066 }
2067 Input:
2068 Reference Text: {reference}
2069 Query: {query}
2070 Claim: {claim}
2071 (Reiteration of the instruction)
2072 Return a JSON5 object with "reasoning" and "answer" (true if plausibly
2073   hallucinated; else false).
2074
2075 Output:

```

2077 A.5.8 MODEL CONFIDENCE SCORE

```

2079 You are an AI assistant tasked with assigning a confidence score to a
2080 claim based on its factuality.
2081
2082 Instructions:
2083 1. You are given:
2084   - A reference text for the query.
2085   - A query.
2086   - A claim made in response to the query.
2087 2. Rate the factuality of the claim with a numeric score in [0.0, 1.0]:
2088   - 0.0 -> The claim is false, contradicts the reference, or is not well
2089   -supported.
2090   - 1.0 -> The claim is factual and well-supported by the reference, or
2091   -can be logically deduced from the reference and query.
2092 3. Provide:
2093   - A score.
2094   - The parts of the reference text that directly support your scoring
2095   -decision.
2096   - A reasoning statement describing your rationale.
2097 4. You must assign a numeric score. Never return null, None, or a non-
2098   -numeric value.
2099
2100 Important: A claim should be considered true (score 1.0) if it is either:
2101 - Directly stated in the reference text, OR
2102 - Can be logically deduced or calculated from the information provided in
2103   -the reference text and query.
2104 For mathematical claims, perform the necessary calculations based on the
2105   -given data.
2106 Output Requirements:
2107 - Output ONLY a single VALID JSON5 object with EXACTLY these keys:
2108   {
2109     "highlighted_text": "Part(s) of the reference text that support the
2110     decision.",
2111     "reasoning": "A reasoning statement describing your rationale.",

```

```

2106     "score": 0.0-1.0
2107   }
2108
2109 JSON5 Rules:
2110 - Use DOUBLE QUOTES (") for all keys and all string values.
2111 - Escape double quotes inside string values as \".
2112 - Escape backslashes as \\.
2113 - No trailing commas in objects or arrays.
2114 - Follow the schema exactly.
2115
2116 Do NOT include:
2117 - Any text, explanations, comments, or formatting outside of the JSON5.
2118
2119 Input:
2120 Reference Text: {reference}
2121 Query: {query}
2122 Claim: {claim}
2123
2124 Reiteration of Instructions:
2125 Return a single JSON5 object with "highlighted_text", "reasoning", "score
2126 ". Assign a numeric score never null/None.
2127
2128 Output:

```

2127 A.5.9 MERGER (FACTSCORE)

```

2129 You will get an instruction and a set of facts that are true. Construct
2130 an answer using ONLY the facts provided, and try to use all facts as
2131 long as its possible. If the input facts are empty, output the empty
2132 string. Do not repeat the instruction.
2133
2134 Input:
2135 The facts: {claims}
2136 The instruction: {query}
2137 Remember, If the input facts are empty, output the empty string. Do not
2138 repeat the instruction.
2139
2140 Output:

```

2139 A.5.10 MERGER (NATURAL QUESTIONS)

```

2141 You will get a natural question and parts of an answer, which you are to
2142 merge into coherent prose. Make sure to include all the parts in the
2143 answer. There may be parts that are seemingly unrelated to the others
2144 , but DO NOT add additional information or reasoning to merge them.
2145 If the input parts are empty, output the empty string. Do not repeat
2146 the question.
2147
2148 Input:
2149 The parts: {claims}
2150 The question: {query}
2151 Remember, DO NOT add any additional information or commentary, just
2152 combine the parts. If the input parts are empty, output the empty
2153 string. Do not repeat the question.
2154
2155 Output:

```

2154 A.5.11 MERGER (MATH)

```

2156 You will get a math problem and a set of steps that are true. Construct
2157 an answer using ONLY the steps provided. Make sure to include all the
2158 steps in the answer, and do not add any additional steps or
2159 reasoning. These steps may not fully solve the problem, but merging
2160 them could assist someone in solving the problem. If the input steps
2161 are empty, output the empty string. Do not repeat the math problem.

```

2160 Input:
 2161 The steps: {claims}
 2162 The math problem: {query}
 2163 Remember, do not do any additional reasoning, just combine the given
 2164 steps. If the input steps are empty, output the empty string. Do not
 2165 repeat the math problem.
 2166 Output:

2167 2168 2169 A.5.12 CORRECTNESS

2171 I need your help in evaluating an answer provided by an LLM against
 2172 ground truth answers. Your task is to determine if the LLMs response
 2173 matches the ground truth answers. Please analyze the provided data
 2174 and make a decision.

2175 Instructions:
 2176 1. Carefully compare the "Predicted Answer" with the "Ground Truth
 2177 Answers".
 2178 2. Consider the substance of the answers look for equivalent information
 2179 or correct answers. Do not focus on exact wording unless the exact
 2180 wording is crucial to the meaning.
 2181 3. Your final decision should be based on whether the meaning and the
 2182 vital facts of the "Ground Truth Answers" are present in the "
 2183 Predicted Answer."
 2184 4. Categorize the answer as one of the following:
 2185 - "perfect": The answer is completely correct and matches the ground
 2186 truth.
 2187 - "acceptable": The answer is partially correct or contains the main idea
 2188 of the ground truth.
 2189 - "incorrect": The answer is wrong or contradicts the ground truth.
 2190 - "missing": The answer is "I dont know", "invalid question", or similar
 2191 responses indicating lack of knowledge.

2190 Output Requirements:
 2191 - Output ONLY a single VALID JSON5 object with EXACTLY this schema:
 2192 {
 2193 "reasoning": "A reasoning statement describing your rationale.",
 2194 "answer": "One of perfect, acceptable, incorrect, or missing"
 2195 }

2196 JSON5 Rules:
 2197 - Use DOUBLE QUOTES (") for all keys and all string values.
 2198 - Escape double quotes inside string values as \".
 2199 - Escape backslashes as \\
 2200 - No trailing commas in objects or arrays.
 2201 - Use the exact top-level container specified and close it properly.
 2202 - Do not include comments, code fences, or any text outside the JSON5
 2203 output.
 2204 - Follow the schema exactly; do not add or omit keys.

2205 Do NOT include:
 2206 - Any text, explanations, comments, or formatting outside of the JSON5.
 2207 - Any code block delimiters (e.g., ``json`).

2208 Input:
 2209 Query: {query}
 2210 Predicted Answer: {merged_string}
 2211 Ground Truth Answer: {answer}

2212 Reiteration of Instructions (before output):
 2213 1. Carefully compare the "Predicted Answer" with the "Ground Truth
 Answers".

```
2214 2. Consider the substance of the answers look for equivalent information
2215 or correct answers. Do not focus on exact wording unless the exact
2216 wording is crucial to the meaning.
2217 3. Your final decision should be based on whether the meaning and the
2218 vital facts of the "Ground Truth Answers" are present in the "
2219 Predicted Answer."
2220 4. Categorize the answer as one of the following:
2221 - "perfect": The answer is completely correct and matches the ground
2222 truth.
2223 - "acceptable": The answer is partially correct or contains the main idea
2224 of the ground truth.
2225 - "incorrect": The answer is wrong or contradicts the ground truth.
2226 - "missing": The answer is "I dont know", "invalid question", or similar
2227 responses indicating lack of knowledge.
2228
2229 Output:
```

2229 A.5.13 SUFFICIENT CORRECTNESS

```
2231 You are an expert LLM evaluator that excels at evaluating a RESPONSE with
2232 respect to a QUERY given a REFERENCE.
2233
2234 Consider the following criteria:
2235 Sufficient Correctness:
2236 1 IF the RESPONSE contains a sufficient amount of CORRECT information (
2237 verified against the REFERENCE) to infer the answer to the QUERY.
2238 0 IF the RESPONSE does not contain a sufficient amount of CORRECT
2239 information to infer the answer to the QUERY.
2240 Important:
2241 Judge only the correctness of the RESPONSE content according to the
2242 REFERENCE.
2243 Do not infer correctness from external knowledge.
2244
2245 Output ONLY a single VALID JSON5 object with EXACTLY these keys:
2246 {
2247   "explanation": An explanation describing your rationale, with with you
2248   will make your decision on sufficient_correctness,
2249   "sufficient_correctness": 1 IF the RESPONSE contains a sufficient
2250   amount of correct information (verified against the REFERENCE) to
2251   infer the answer to the QUERY, 0 IF the RESPONSE does not contain a
2252   sufficient amount of correct information to infer the answer to the
2253   QUERY.
2254 }
2255 - "sufficient_correctness" must be an integer (0/1).
2256
2257 JSON5 Rules:
2258 - Use DOUBLE QUOTES (") for all keys and all string values.
2259 - Escape double quotes inside string values as \".
2260 - Escape backslashes as \\.
2261 - No trailing commas in objects or arrays.
2262 - Use the exact top-level container specified and close it properly.
2263 - Do not include comments, code fences, or any text outside the JSON5
2264 output.
2265 - Follow the schema exactly; do not add or omit keys.
2266
2267 Do NOT include:
2268 - Any text, explanations, comments, or formatting outside of the JSON5.
2269 - Any code block delimiters (e.g., ```json).
2270
2271 Input:
2272 ### QUERY
2273 {query}
2274 ### REFERENCE
2275 {reference}
```

2268 ### RESPONSE
2269 {response}
2270 Output:

2273 A.5.14 MATH REFERENCE

2276 You are a helpful assistant that extracts the prerequisite mathematics
2277 knowledge needed to answer a given question without solving it.

2278 Instructions:

- 2279 1. Read the question.
- 2280 2. Identify only the minimal prerequisite items across concepts,
2281 definitions, theorems/properties, formulas, techniques, notation,
2282 assumptions/conditions, and common pitfalls required to answer the
2283 question.
- 2284 3. When you state a theorem, a definition, or a formula, write out its
2285 full, standard statement (not just the name). For theorems, include
2286 hypotheses and conclusions; for definitions, give the precise meaning
2287 ; for formulas, write the exact equation(s) in standard notation.
- 2288 4. Do NOT provide the answer or partial solution steps.
- 2289 5. Do NOT use examples that arise directly from the given question; keep
statements general and problem-agnostic.

2290 Output Requirements:

- 2291 - Produce PLAIN TEXT ONLY.
- 2292 - Write free-form prose (no headings, lists, or numbering).
- 2293 - Mention required topics, definitions, fully written theorems/properties
2294 , formulas, techniques, notation conventions, assumptions/conditions,
and common pitfalls in sentences.
- 2295 - The output may be long; include complete statements where needed.

2296 Do NOT include:

- 2297 - Any answer, hints, or step-by-step solution.
- 2298 - Any examples derived from the given question.

2300 Reiteration of the instructions:

- 2301 1. Read the question.
- 2302 2. Identify only the minimal prerequisite items across concepts,
2303 definitions, theorems/properties, formulas, techniques, notation,
2304 assumptions/conditions, and common pitfalls required to answer the
question.
- 2305 3. When you state a theorem, a definition, or a formula, write out its
2306 full, standard statement (not just the name). For theorems, include
2307 hypotheses and conclusions; for definitions, give the precise meaning
2308 ; for formulas, write the exact equation(s) in standard notation.
- 2309 4. Do NOT provide the answer or partial solution steps.
- 2310 5. Do NOT use examples that arise directly from the given question; keep
2311 statements general and problem-agnostic.

2312 Input:

2313 {query}

2314 Output:

2318 B ETHICS STATEMENTS

2320 This work adheres to the ICLR Code of Ethics. We do not anticipate any direct risks of harm to
2321 individuals or groups arising from this research.

2322 C REPRODUCIBILITY STATEMENT
2323

2324 We are committed to ensuring the reproducibility and transparency of our work. To this end, we
2325 will fully open-source our code upon publication. The released repository will include the complete
2326 implementation scripts. All datasets and models used in this study are already publicly available.
2327

2328 D USAGE OF LARGE LANGUAGE MODELS
2329

2330 Large language models were employed during the preparation of this manuscript. Commercial
2331 LLMs were used solely for editing purposes, improving clarity, grammar, and adherence to academic
2332 style, without altering the authors' intended meaning or substantive contributions.
2333

2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375