

# Flaw or Artifact? Rethinking Prompt Sensitivity in Evaluating LLMs

Andong Hua<sup>1\*</sup>, Kenan Tang<sup>1\*</sup>, Chenhe Gu<sup>2</sup>, Jindong Gu<sup>3</sup>, Eric Wong<sup>4</sup>, Yao Qin<sup>1</sup>

<sup>1</sup> UC Santa Barbara, <sup>2</sup> UC Irvine, <sup>3</sup> University of Oxford, <sup>4</sup> University of Pennsylvania  
dongx1997@ucsb.edu, yaoqin@ucsb.edu

## Abstract

Prompt sensitivity, referring to the phenomenon where paraphrasing (i.e., repeating something written or spoken using different words) leads to significant changes in large language model (LLM) performance, has been widely accepted as a core limitation of LLMs. In this work, we revisit this issue and ask: Is the widely reported high prompt sensitivity truly an inherent weakness of LLMs, or is it largely an artifact of evaluation processes? To answer this question, we systematically evaluate 7 LLMs (e.g., GPT and Gemini family) across 6 benchmarks, including both multiple-choice and open-ended tasks on 12 diverse prompt templates. We find that much of the prompt sensitivity stems from heuristic evaluation methods, including log-likelihood scoring and rigid answer matching, which often overlook semantically correct responses expressed through alternative phrasings, such as synonyms or paraphrases. When we adopt LLM-as-a-Judge evaluations, we observe a substantial reduction in performance variance and a consistently higher correlation in model rankings across prompts. Our findings suggest that modern LLMs are more robust to prompt templates than previously believed, and that prompt sensitivity may be more an artifact of evaluation than a flaw in the models.

## 1 Introduction

Large Language Models (LLMs) have achieved remarkable success across a wide range of tasks (Hua et al., 2024; Liévin et al., 2024). Moreover, LLMs are good at following diverse instructions, so that users are not required to follow a fixed template when asking a question. This has led to concerns about prompt sensitivity, where differences in prompt phrasing can substantially affect benchmark performance, casting doubt on the reliability of evaluations (Polo et al., 2024; Mizrahi et al.,

2024; Chatterjee et al., 2024; Sclar et al., 2023). More critically, the relative rankings of LLMs can shift substantially depending on the prompt template used (Polo et al., 2024; Mizrahi et al., 2024). For example, simply changing the option format from letters (e.g., “A:”) to numbers (e.g., “(1)”) completely *reverses* the ranking order of four evaluated open-source models in ARC-Challenge (Clark et al., 2018).

Although existing studies have reported that LLMs are highly sensitive to prompt phrasing (Voronov et al., 2024; Mizrahi et al., 2024), this remains counterintuitive given that instruction-tuned LLMs are explicitly optimized to handle a wide range of input formats. For example, instruction-tuning datasets such as FLAN (Longpre et al., 2023) and Super-NaturalInstructions (Wang et al., 2022) include a diverse collection of tasks (e.g., question answering, summarization, classification) with varying natural language prompt templates (Zhang et al., 2023). This contradiction raises a critical question:

*Is prompt sensitivity an inherent flaw in LLMs, or merely an artifact of the evaluation process?*

To investigate this, we find that previous studies (Voronov et al., 2024; Chatterjee et al., 2024) typically rely on heuristic evaluation, such as regular-expression-based answer extraction or log-likelihood scoring over candidates. These heuristic evaluation approaches, though historically popular due to their simplicity (Zellers et al., 2019; Reddy et al., 2019; Brown et al., 2020; Hendrycks et al., 2021a), may introduce errors when model outputs deviate from expected formats. More specifically, models may generate correct answers, but because their outputs are not aligned with the rigid evaluation format, they are mistakenly marked as incorrect. This issue becomes more pronounced as models have become more open-ended and diverse in their output formats (Hurst et al., 2024; Yang et al., 2025), potentially leading to inflated esti-

\*Equal contributions.

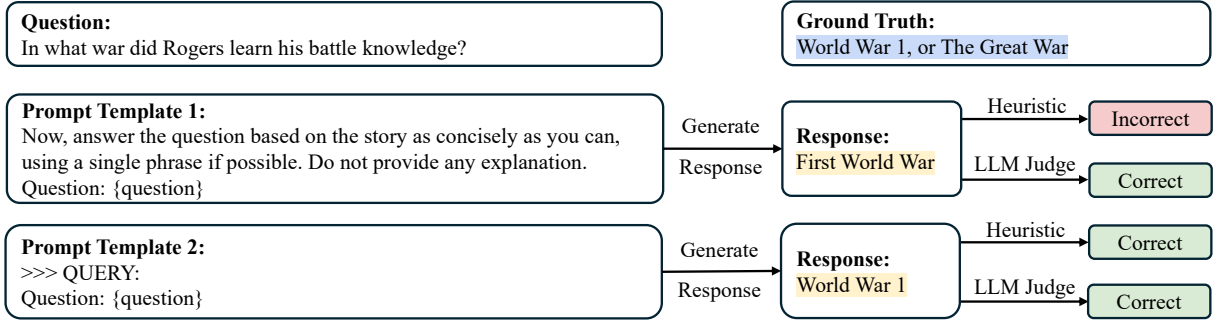


Figure 1: When provided with diverse prompt templates, LLMs provide different but semantically equivalent responses. Heuristic evaluation fails to match the different answers with the ground truth, exaggerating prompt sensitivity. In contrast, an LLM judge is able to identify the semantic equivalence consistently.<sup>1</sup>

mates of prompt sensitivity (Figure 1).

To rigorously assess whether LLMs truly suffer from prompt sensitivity, we revisit this issue using a more robust evaluation strategy: LLMs as judges. Now widely adopted in recent benchmarks (Wei et al., 2024; Zheng et al., 2023), this approach shifts evaluation from rigid pattern-matching to semantic assessment, enabling a more reliable examination of prompt sensitivity. Compared to heuristics, LLM judges can better handle various output formats, paraphrasing, and ambiguous cases, making them more aligned with human evaluation (Liu et al., 2023; Kocmi and Federmann, 2023; Zheng et al., 2023; Wang et al., 2023).

Using both heuristic methods and LLM-as-a-Judge, we conduct a comprehensive evaluation of four open-source and three closed-source LLMs across 12 prompt templates (without cherry-picking) and six diverse benchmarks, including both multiple-choice and open-ended generation tasks. We discover that heuristic evaluation methods often **exaggerate** the prompt sensitivity of LLMs. For instance, on ARC-Challenge, the performance of Gemma-2.0 varies widely across prompts, with accuracy ranging from 0.25 to 0.90 and a high standard deviation of 0.28. In contrast, when evaluated using LLM-as-a-Judge, its accuracy varies by only 0.17 across the same set of prompt templates, with a much lower standard deviation of just 0.005. Furthermore, the *Spearman rank correlation* measuring model performance rankings across the four open-source models remarkably improves from 0.31 under heuristic evaluation to 0.92 with LLM-as-a-Judge. These findings suggest that previously reported prompt sensitivity may be significantly overstated due to the limita-

tions of heuristic evaluation. When evaluated with LLM-as-a-Judge, models show far more consistent performance and stable rankings across prompts. Notably, we have conducted a comprehensive human study to verify the reliability of LLM-judges. The high consistency between LLM-judges and human annotators strongly indicates that prompt sensitivity is largely an artifact of the evaluation method rather than an inherent flaw in LLMs.

## 2 Method

Our method consists of three main steps: diverse prompt template construction (Section 2.1), LLM-as-a-Judge evaluation (Section 2.2), and prompt template sensitivity measurement (Section 2.3).

### 2.1 Diverse Prompt Template Construction

To evaluate LLM sensitivity to prompt phrasing, we construct diverse prompt templates for each benchmark. These templates vary in instruction wording, answer formatting (e.g., using letters vs. numbers), and how responses are requested, while keeping the task content unchanged (Appendix A).

In practice, we use GPT-4o to paraphrase the original prompts. For multiple-choice datasets, we create a shared pool of 12 diverse templates used across all benchmarks. For open-ended generation tasks, we generate 12 templates per benchmark to better accommodate domain-specific styles.

### 2.2 LLM-as-a-Judge Evaluation

Heuristic evaluation methods often fail when model outputs deviate from expected formats. To address this limitation, we adopt LLMs as robust judges. In this approach, an LLM judge is given the original question, the correct answer, and the model’s predicted response. The judge is prompted to determine whether the predicted response semantically

<sup>1</sup>This is an example from NarrativeQA. Heuristic method uses word-level F1; “Incorrect” is shown here for illustration purposes, indicating a lower score than a correct answer.

matches the correct answer (Appendix B).

While the overall format remains consistent, we introduce minor benchmark-specific adjustments to the judging prompt. For example, for GPQA, the judge is instructed to “Ignore all explanation” to ensure it focuses solely on the final answer.

### 2.3 Prompt Sensitivity Measurement

We measure sensitivity with two metrics: *performance variation* and *ranking consistency*.

**Performance Variation.** For each model and dataset, we compute the accuracy under every prompt template and report the standard deviation across all prompt variants. Let  $P = \{p_1, p_2, \dots, p_n\}$  denote the set of prompt templates for a given benchmark  $D$ , and let  $f$  be the model under evaluation. The performance of model  $f$  under prompt  $p_i$  is denoted by  $A_{f,D}^{p_i}$ . The prompt sensitivity of model  $f$  on dataset  $D$  is then quantified as:  $\text{std}_f = \text{StdDev} \left( \left\{ A_{f,D}^{p_i} \right\}_{i=1}^n \right)$ . A lower standard deviation indicates that the model’s performance is stable across different prompt templates.

**Ranking Consistency.** Beyond absolute performance, we also measure how model rankings vary across prompt templates. Given a set of  $K$  models, we rank them based on their performance under each prompt and compute pairwise *Spearman’s rank correlation* between all prompt pairs (Spearman, 1904). Given two templates  $p_i$  and  $p_j$ , and the corresponding performance vectors  $\{A_{f_k,D}^{p_i}\}_{k=1}^K$  and  $\{A_{f_k,D}^{p_j}\}_{k=1}^K$ , we calculate Spearman’s rank correlation coefficient  $\rho_{ij} = 1 - \frac{6 \sum_{k=1}^K d_k^2}{K(K^2-1)}$ , where  $d_k$  is the difference in rankings of the  $k$ -th model under prompts  $p_i$  and  $p_j$ , and  $K$  is the number of models. The rank correlation coefficient,  $\rho$ , ranges from  $-1$  to  $1$ , with higher values indicating stronger agreement in ranking consistency.

To measure overall ranking consistency, we compute the mean Spearman’s rank correlation coefficient, denoted as  $\bar{\rho}$ , across all pairs of prompt templates. This mean score  $\bar{\rho}$  serves as a comprehensive metric for evaluating the stability of model rankings under prompt variation. A higher  $\bar{\rho}$  suggests that evaluations are more robust and less dependent on the specific prompt phrasing.

## 3 Results and Discussion

### 3.1 Experimental Setup

**Models.** We evaluate LLaMA-3.1-8B-Instruct (LLaMA-3.1) (Dubey et al., 2024), Qwen2-7B-

Instruct (Qwen-2) (Yang et al., 2024), Gemma-2-9B-it (Gemma-2) (Gemma Team et al., 2024), Ministral-8B-Instruct-v0.2 (Ministral) (MistralAI, 2024), GPT-4o-mini (July 2024), GPT-4.1-mini (April 2025), and Gemini 2.0 Flash (February 2025).

**Benchmarks.** We evaluate on six benchmarks covering both multiple-choice and open-ended tasks. The multiple-choice datasets include ARC-Challenge (Clark et al., 2018), GPQA-diamond (Rein et al., 2024), and OpenbookQA (Mihaylov et al., 2018), where answers are selected from discrete options (e.g., A/B/C/D). For these tasks, heuristic evaluation uses log-likelihood scoring over answer options. The open-ended datasets include NarrativeQA (Kočíský et al., 2018), MATH (Hendrycks et al., 2021b), and SimpleQA (Wei et al., 2024), where model responses are free-form. For NarrativeQA and MATH, heuristic evaluation applies format-specific extraction and normalization (see Appendix C). For SimpleQA, no rule-based parser is available, so we report results only under the LLM-as-a-Judge framework. All evaluations use greedy decoding to ensure deterministic outputs.

**LLM-as-a-Judge.** Unless otherwise specified, we employ Gemini 2.0 Flash as the LLM-as-a-Judge across all benchmarks.

| Dataset        | $\bar{\rho}_{\text{Heuristic}}$ | $\bar{\rho}_{\text{LLM}}$ |
|----------------|---------------------------------|---------------------------|
| ARC-Challenge* | 0.3036                          | 0.9546 (0.9187)           |
| OpenbookQA*    | 0.4212                          | 0.9386 (0.7360)           |
| GPQA Diamond*  | 0.1542                          | 0.8960 (0.5048)           |
| NarrativeQA†   | 0.5927                          | 0.8662                    |
| MATH           | 0.9593                          | 0.9647                    |
| SimpleQA       | –                               | 0.8121                    |

Table 1: Average Spearman rank correlation ( $\bar{\rho}$ ) across prompt templates using heuristic vs. LLM-as-a-Judge. Results use all 7 models unless noted: \* heuristic uses only 4 open-source models (parentheses show LLM-as-a-Judge restricted to the same 4 models); † NarrativeQA uses LLaMA-3.1 and 3 proprietary models due to context limits; “–” indicates heuristic not applicable.

### 3.2 Heuristic Evaluation Exaggerates Prompt Sensitivity of LLMs

When comparing the performance of the model under heuristic evaluation and LLM-as-a-Judge, we find that the heuristic methods exhibit significantly greater sensitivity to prompt variation (Figure 2). On ARC-Challenge, all open-source models except

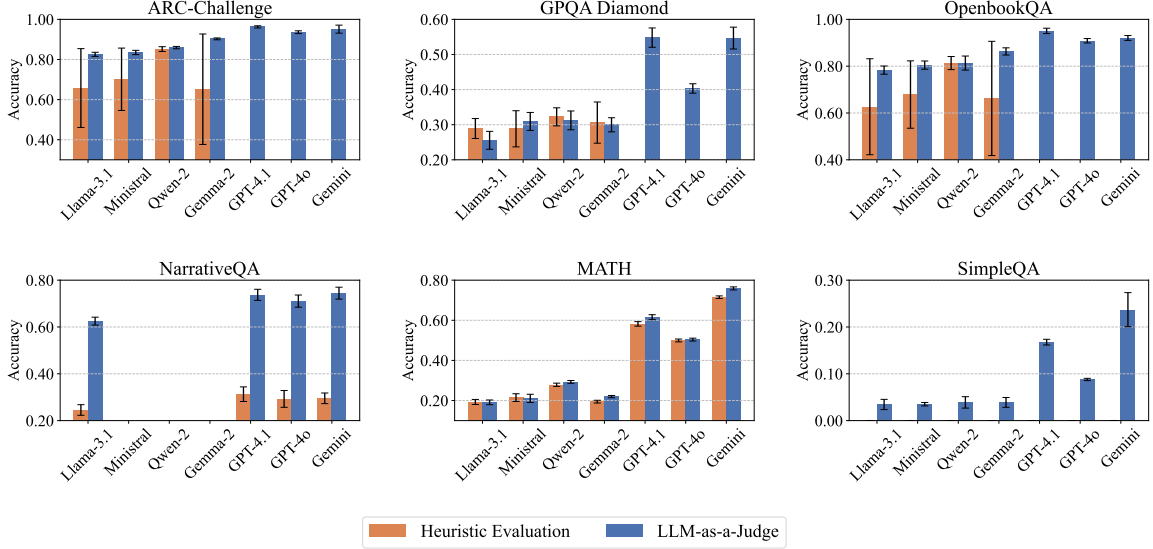


Figure 2: The mean and standard deviation of performance across different prompt templates. For all 6 datasets, we show the statistics for all pairs of evaluation methods and models, excluding the cases when the model’s context length is not enough for the task or when the heuristic evaluation method is not available. The standard deviation of the LLM-as-a-Judge method is always low. For NarrativeQA, the absence of results for Mistral and Qwen is primarily due to the long-context requirement of this dataset.

Qwen-2 show much higher standard deviations under heuristics. For instance, Gemma-2.0 yields a deviation of 0.28, versus just 0.005 with LLM-as-a-Judge. Its accuracy range spans 0.25–0.90 under heuristics, compared to only 0.17 with LLM-as-a-Judge. Additionally, mean accuracy improves under LLM-as-a-Judge, suggesting heuristic methods often miss valid answers due to overly rigid extraction rules.

Beyond variance in accuracy, we also assess ranking consistency across prompts (Section 2.3). On ARC-Challenge, the average Spearman rank correlation across prompts among open-source models increases from 0.30 (heuristics) to 0.92 (LLM-as-a-Judge), and further to 0.95 when proprietary models are included. On NarrativeQA, the correlation rises from 0.40 (heuristics) to 0.87 (LLM-as-a-Judge). These findings suggest that prompt sensitivity observed in prior work is largely an artifact of heuristic evaluation, not an inherent flaw of LLMs.

**Well-designed heuristic methods show low prompt sensitivity similar to LLM-as-a-Judge.** For MATH (Hendrycks et al., 2021b), the heuristic approach incorporates symbolic simplification, expression normalization, and equivalence checking using tools such as sympy. Under these conditions, we observe prompt sensitivity results that are comparable to those obtained via LLM-as-a-Judge evaluation, with similarly low accuracy variance

and high ranking consistency. These results indicate that, with sufficient domain-specific prompt engineering, heuristic methods can provide stable evaluations. This further supports the conclusion that modern LLMs exhibit less prompt-induced variance than previously reported.

**Newly proposed benchmarks exhibit low prompt sensitivity.** We extend the evaluation to SimpleQA (Wei et al., 2024), a newly proposed benchmark for factual and commonsense reasoning. As no official heuristic evaluation is available, we use LLM-as-a-Judge by default. Applying our prompt sensitivity analysis, we observe a low standard deviation and a high Spearman rank correlation of 0.8121 across 12 prompt templates. These results indicate that model performance remains stable across different prompt variations.

### 3.3 Consistency of LLM-as-a-Judge Evaluation Across LLM Judges

To further examine whether evaluation outcomes vary significantly when using different LLMs as judges, we conducted an additional analysis on the ARC-Challenge dataset with GPT-4o-mini as the judge. As shown in Table 2, both the standard deviation across different prompts and the ranking correlation remain consistent regardless of whether GPT or Gemini is used as the judge. This suggests that LLM-as-a-Judge evaluations are robust across different LLMs, reinforcing our main argument that



| Model            | GPT Judge           | Gemini Judge        |
|------------------|---------------------|---------------------|
| Llama-3.1        | 0.8276 $\pm$ 0.0094 | 0.8257 $\pm$ 0.0102 |
| Mistral          | 0.8349 $\pm$ 0.0105 | 0.8349 $\pm$ 0.0107 |
| Qwen2            | 0.8590 $\pm$ 0.0063 | 0.8592 $\pm$ 0.0059 |
| Gemma-2          | 0.9027 $\pm$ 0.0048 | 0.9027 $\pm$ 0.0046 |
| Gemini           | 0.9590 $\pm$ 0.0050 | 0.9512 $\pm$ 0.0200 |
| GPT-4.1          | 0.9617 $\pm$ 0.0037 | 0.9626 $\pm$ 0.0054 |
| GPT-4o           | 0.9371 $\pm$ 0.0050 | 0.9360 $\pm$ 0.0069 |
| Rank Correlation | 0.9621              | 0.9546              |

Table 2: Accuracy scores (Mean  $\pm$  Std) across prompts in ARC-Challenge with different LLM judges. Results show high consistency between GPT and Gemini judges.

prompt sensitivity is largely a byproduct of heuristic evaluation rather than true model instability.

#### 4 LLM-as-a-Judge Evaluation Aligns with Human Annotations

To assess the reliability of LLM-judges, we compare them against human annotations. We recruit human annotators to manually evaluate answer correctness. For each dataset, we randomly sample 50 questions and collect answers from one model under all 12 prompt templates, yielding 600 answers per dataset-model pair. For ARC-Challenge, OpenbookQA, and GPQA Diamond, we evaluate Gemma-2; for NarrativeQA, MATH, and SimpleQA, we evaluate GPT-4.1-mini. We also report results on the combined set of all six subsets.

Human annotators evaluate whether each answer matches the corresponding ground truth. Their judgments are highly consistent, and majority voting is used to resolve the few cases of disagreement. Moreover, the majority-voted outcomes closely align with the LLM-judge results, underscoring the reliability of LLM-as-a-Judge evaluation. More details and the original human annotation instructions can be found in Appendix D.

From the results, we draw two key observations.

**Observation 1** Human-annotated results show consistently high agreement on the answer correctness (Table 3), with a high Fleiss’  $\kappa$  over 0.6 (Fleiss, 1971). This provides strong evidence that the answer correctness does not vary substantially across different prompt templates.

**Observation 2** The rate of perfect agreement is also notably high (Table 3). We define perfect agreement as cases where the answer correctness is identical across all prompt templates for a given question. This further reinforces our conclusion

that prompt templates have minimal impact on the answer correctness.

| Dataset       | Agreement | Perfect Agreement |
|---------------|-----------|-------------------|
| Arc-Challenge | 0.7654    | 86%               |
| OpenbookQA    | 0.7250    | 80%               |
| GPQA-Diamond  | 0.6708    | 52%               |
| NarrativeQA   | 0.7258    | 66%               |
| MATH          | 0.7859    | 68%               |
| SimpleQA      | 0.7881    | 88%               |
| Combined      | 0.8132    | 73%               |

Table 3: Correctness of answers shows minimal variation across prompt templates, based on human annotations. “Agreement” is measured by Fleiss’  $\kappa$  (Fleiss, 1971) across 12 prompt templates (higher is better). “Perfect Agreement” is defined as cases where all 12 prompt templates yield the same correctness judgment for a given question. We report the percentage of questions with perfect agreement.

#### 5 Related Work

The ranking inconsistency with diverse prompt templates has been widely reported (Polo et al., 2024; Mizrahi et al., 2024; Chatterjee et al., 2024; Sclar et al., 2023). Mizrahi et al. (2024) conducted a large-scale study showing significant accuracy differences across prompt variants. Voronov et al. (2024) further showed that no prompt format consistently performs best across models. To address this, prior work often assumes that LLMs are inherently unstable to prompt changes. For example, Polo et al. (2024) estimates the distribution of accuracy across prompts to improve evaluation efficiency. However, all existing methods attribute the sensitivity to model behavior. In contrast, we show that a key factor is the heuristic evaluation protocol itself, which often leads to misclassification of correct outputs and overstates prompt sensitivity.

#### 6 Conclusion

In this work, we demonstrate that much of the observed prompt sensitivity in LLM evaluations is not due to inherent model weaknesses, but rather an artifact introduced by heuristic evaluation methods. Through comprehensive experiments using LLM-as-a-Judge across multiple benchmarks and prompt templates, we reveal that model performance and rankings are substantially more stable and reliable than previously reported. We hope this work sheds light on prompt sensitivity in LLM evaluation and encourages broader adoption of LLM-as-a-Judge to evaluate the true capabilities of LLMs.

## Limitations

Due to computational constraints, we evaluate each benchmark using only 12 prompt templates. However, we find that results are stable across scales. For example, on ARC-Challenge, the ranking consistency and variance metrics using 12 prompts closely match those obtained with over 100 prompts, suggesting that our analysis is representative.

## References

- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. [LongBench: A bilingual, multi-task benchmark for long context understanding](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137, Bangkok, Thailand. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Anwoy Chatterjee, HSVNS Kowndinya Renduchintala, Sumit Bhatia, and Tanmoy Chakraborty. 2024. Posix: A prompt sensitivity index for large language models. *arXiv preprint arXiv:2410.02185*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Andong Hua, Mehak Preet Dhaliwal, Ryan Burke, Laya Pullela, and Yao Qin. 2024. Nutribench: A dataset for evaluating large language models on nutrition estimation from meal descriptions. *arXiv preprint arXiv:2407.12843*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The narrativeqa reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6:317–328.
- Tom Kocmi and Christian Federmann. 2023. [Large language models are state-of-the-art evaluators of translation quality](#). In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.
- Victor Liévin, Christoffer Egeberg Hother, Anna Guldburger Motzfeldt, and Ole Winther. 2024. [Can large language models reason about medical questions?](#) *Patterns*, 5(3):100943.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, and 1 others. 2023. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pages 22631–22648. PMLR.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*.
- MistralAI. 2024. [Ministral-8b-instruct-2410. https://huggingface.co/mistralai/Ministral-8B-Instruct-2410](https://huggingface.co/mistralai/Ministral-8B-Instruct-2410).

- Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2024. State of what art? a call for multi-prompt llm evaluation. *Transactions of the Association for Computational Linguistics*, 12:933–949.
- Felipe Maia Polo, Ronald Xu, Lucas Weber, Mírian Silva, Onkar Bhardwaj, Leshem Choshen, Allysson Flavio Melo de Oliveira, Yuekai Sun, and Mikhail Yurochkin. 2024. Efficient multi-prompt evaluation of llms. *arXiv preprint arXiv:2405.17202*.
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. [CoQA: A conversational question answering challenge](#). *Transactions of the Association for Computational Linguistics*, 7:249–266.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. [GPQA: A graduate-level google-proof q&a benchmark](#). In *First Conference on Language Modeling*.
- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.
- Charles Spearman. 1904. [The proof and measurement of association between two things](#). *The American Journal of Psychology*, 15(1):72–101. JSTOR 1412159.
- Anton Voronov, Lena Wolf, and Max Ryabinin. 2024. Mind your format: Towards consistent evaluation of in-context learning improvements. *arXiv preprint arXiv:2401.06766*.
- Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. [Is ChatGPT a good NLG evaluator? a preliminary study](#). In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11, Singapore. Association for Computational Linguistics.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, and 1 others. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109.
- Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. 2024. Measuring short-form factuality in large language models. *arXiv preprint arXiv:2411.04368*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. [Qwen3 technical report](#). *arXiv preprint arXiv:2505.09388*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. [Qwen2 technical report](#). *Preprint, arXiv:2407.10671*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and 1 others. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.

## A Diverse Prompts

In this section, we list the diverse prompt templates we use for each benchmark.

For ARC-Challenge, GPQA, and OpenbookQA, we use the following 12 prompt templates:

1. Evaluate the choices and select the most appropriate answer.\n{question}\n\nThe options are as follows:\nOption A: {first\_option}\nOption B: {second\_option}\nOption C: {third\_option}\nOption D: {fourth\_option}\n\nYour answer should be formatted as:\n'I have chosen option [choice]'\n\nwhere [choice] is your selected answer.\n
2. Review the available options and select the one you think is correct.\n{question}\n\nAvailable answers include:\nA ) {first\_option}\nB ) {second\_option}\nC ) {third\_option}\nD ) {fourth\_option}\n\n\nResponse:
3. Select the correct answer based on your understanding.\n{question}\n\nPick from the following options:\n[A] {first\_option}\n[B] {second\_option}\n[C] {third\_option}\n[D] {fourth\_option}\n\nPlease respond with 'Option [choice]'.\n

4. Evaluate the options presented and select the most suitable. {question}\nAvailable answers:\n[A] {first\_option}\n[B] {second\_option}\n[C] {third\_option}\n[D] {fourth\_option}\n\nExpress your choice as: 'The answer is [choice].'\nwhere [choice] is your selected option.\n
  5. Based on the question presented, choose the most fitting response. {question}\nAvailable answers are:\nA: {first\_option}\nB: {second\_option}\nC: {third\_option}\nD: {fourth\_option}\nPlease provide your response in the following format:\n'Your choice: [option]'\nwhere [option] corresponds to the letter or number you selected.\n
  6. From the options below, select the response that you believe is correct. {question}\nChoices to consider:\n1. {first\_option}\n2. {second\_option}\n3. {third\_option}\n4. {fourth\_option}\nResponse:
  7. Select your answer from the provided list of options.\n{question}\nOptions are:\nThe choice is A: {first\_option}\nThe choice is B: {second\_option}\nThe choice is C: {third\_option}\nThe choice is D: {fourth\_option}\n\nChoose your answer:
  8. After considering the options, choose the best possible answer. {question}\n\nThe following choices are available:\nA: {first\_option}\nB: {second\_option}\nC: {third\_option}\nD: {fourth\_option}\n\nState your answer as:\n'Answer: [choice]'\n
  9. Analyze the selections and provide your choice.\n{question}\n\nYour options are listed below:\nOption 1 - {first\_option}\nOption 2 - {second\_option}\nOption 3 - {third\_option}\nOption 4 - {fourth\_option}\n\nYour response:
  10. Consider the following question and determine the right response.\n{question}\n\nWhich of the following answers do you prefer?\nOption 1: {first\_option}\nOption 2: {second\_option}\nOption 3: {third\_option}\nOption 4: {fourth\_option}\n\nI select:
  11. Determine which option best answers the question asked.\n{question}\n\nPossible choices are as follows:\nOption [A] {first\_option}\nOption [B] {second\_option}\nOption [C] {third\_option}\nOption [D] {fourth\_option}\n\nFinal answer:
  12. Identify the option that best answers the question posed. {question}\n\nConsider these choices:\nSelect option 1: {first\_option}\nSelect option 2: {second\_option}\nSelect option 3: {third\_option}\nSelect option 4: {fourth\_option}\n\nChoice provided:
- In the prompt templates, {question} is the question, and {first\_option}, {second\_option}, {third\_option}, and {fourth\_option} are the options.
- For NarrativeQA, we use the following 12 prompt templates:
1. You are given a story, which can be either a novel or a movie script, and a question. Answer the question as concisely as you can, using a single phrase if possible. Do not provide any explanation.\n\nStory: {context}\n\nNow, answer the question based on the story as concisely as you can, using a single phrase if possible. Do not provide any explanation.\n\nQuestion: {question}\n\nAnswer:
  2. Below is an excerpt from a mystery or thriller story, followed by a question. Provide the most accurate answer you can in a single phrase or sentence fragment. No elaboration is needed.\n\nStory: {context}\n\nExamine the situation carefully and respond.\n\nQuestion: {question}\n\nAnswer:
  3. You are presented with a passage from literary fiction or cinematic writing and a comprehension question. Respond succinctly with a phrase. Avoid any additional commentary.\n\nStory: {context}\n\nAnalyze and respond concisely.\n\nQuestion: {question}\n\nAnswer:
  4. A tale from a distant world or magical land is told below, followed by a question from a curious scholar. Give your answer using only



- a few words. No need to explain the lore.  
 Story: {context}  
 What say you?  
 Question: {question}  
 Answer:
5. You're reading a gritty tale from the backstreets of the city. A question follows. Keep your answer clipped, clean, and under the radar—just a phrase, no fluff.  
 Story: {context}  
 Here's the case:  
 Question: {question}  
 Answer:
  6. Welcome to *\*Plot Points\**! We'll give you a story snippet and a question—your job is to give the fastest, most precise answer possible. One phrase, no lifelines!  
 Story: {context}  
 Let's play!  
 Question: {question}  
 Answer:
  7. The record shows the following account. A question will now be entered into the record. Provide your answer in a short, factual phrase. No commentary permitted.  
 Story: {context}  
 Deposition  
 Question:  
 Question:  
 Question:  
 Answer:
  8. Read the excerpt. Answer the question. Keep it short.  
 Story: {context}  
 Question:  
 Answer:
  9. Accessing archive... Story fragment retrieved from Galactic Chronicles. A query follows. Respond with the most relevant concept or phrase. Do not explain.  
 Story: {context}  
 >>> QUERY:  
 Question:  
 Answer:  
 >>> RESPONSE:  
 Answer:
  10. Once upon a time, a story was told. Now a little question is asked. Answer it kindly and briefly—just a few words will do. No need to explain why.  
 Story: {context}  
 Here comes the question:  
 Question:  
 Answer:
  11. From the folds of a lyrical tale, a question emerges like morning light. Respond with a single phrase, a shard of truth—no more, no less.  
 Story: {context}  
 Whisper your reply:  
 Question: {question}  
 Answer:
  12. Intel received. Narrative extracted. Stand by for situational query. Your task: deliver the answer in minimal terms. Do not elaborate.  
 Story: {context}  
 Mission Query:  
 Question: {question}  
 Answer:

In the prompt templates, {context} is the context, and {question} is the question.

For MATH, we use the prompt template {text1}{question}{text2}, where {text1} and {text2} are two strings that enclose the question. The following 12 pairs of ({text1}, {text2}) are used:

1. ({empty\_string}, \nAnswer:\n)
2. (Problem:\n, \nAnswer:\n)
3. (Problem:\n, \nAnswer:\n)
4. (Task:\n\n, \n\nSolution:)
5. (Solve the following math problem:\n\n, \nAnswer:\n)
6. (Solve the following math problem:\n\n, \nAnswer:\n)
7. (\*\*Problem Statement\*\*:\n\n, \n\nSolution:)
8. (Problem:\n, \n\nSolution:)
9. (Solve the following math problem:\n\n, \n\nSolution:)
10. (\*\*Problem Statement\*\*:\n\n, \n\nSolution:)
11. (\*\*Problem Statement\*\*:\n\n, \nAnswer:\n)
12. ({empty\_string}, \nAnswer:\n)

Since MATH uses few-shot prompting for evaluation, we further change the examples provided for each prompt template. Hence, while two pairs of ({text1}, {text2}) could be the same, the actual prompt template is different.

For SimpleQA, we use the prompt template {instruction}{question}, where {question} is the original questions in the benchmark, and {instruction} is one of the 12 following strings:

1. {empty\_string}
2. Ready your reasoning—consider the challenge that follows.  
 \n\n
3. Take a thoughtful pause, then craft your best response to the prompt beneath this line.  
 \n\n
4. Showcase your insight by addressing the upcoming question.  
 \n\n
5. Put your analytical lens on and dive into the inquiry below.  
 \n\n

You are an AI assistant that determines whether a model's prediction matches a given reference answer for a question.  
 You will be given:  
 - A question  
 - A reference (correct) answer  
 - A model's predicted answer  
 Your task is to judge whether the prediction matches the reference answer.  
 Ignore any explanation in the prediction—only the final selected answer matters.  
 Respond with a JSON object in the following format:  
 {'match': true, 'reason': '...'} if the prediction matches the reference  
 {'match': false, 'reason': '...'} if it does not match  
 Be specific and concise in your reasoning.  
 \*\*Do not answer the question or provide any other information.\*\*  
 Here is the input:  
 Question: {question}  
 Reference Answer: {reference\_answer}  
 Model Prediction: {model\_prediction}

Figure 3: An example of a judging prompt. After filling in the question, reference answer, and model prediction, we send the prompt to an LLM judge to get the result.

6. Channel your inner detective: examine the next question and present your findings.\n\n
7. Let your knowledge shine—respond thoughtfully to the statement that follows.\n\n
8. Engage your critical thinking skills and tackle the question that appears next.\n\n
9. Apply the concepts you’ve mastered to answer the forthcoming inquiry.\n\n
10. Use evidence and reasoning to construct your answer to the question below.\n\n
11. Approach the next problem with curiosity and craft a clear solution.\n\n
12. Demonstrate what you’ve learned by addressing the prompt that follows.\n\n

## B Prompts for LLM-as-a-Judge

Figure 3 shows the prompt we use for LLM-as-a-Judge. For each benchmark, we make minor task-specific modifications to the judging prompt. For SimpleQA, we use the official prompt from OpenAI.

## C Heuristic Evaluation Details

**NarrativeQA.** Following LongBench (Bai et al., 2024), we compute word-level F1 overlap between normalized predictions and references. Normalization includes lowercasing, removing punctuation and articles (a, an, the), and collapsing whitespace. For example, the prediction “Fifty years” and the

reference “50 years” would result in a partial overlap and an F1 score of 0.5.

**MATH.** Heuristic evaluation typically extracts the final answer from LaTeX-formatted expressions such as `\boxed{...}` or `\fbox{...}`. This method assumes the model explicitly marks its answer with these delimiters. We follow this approach in our experiments. See Appendix A for examples.

**SimpleQA.** No official heuristic parser is currently available, making rule-based evaluation infeasible. We therefore rely solely on the LLM-as-a-Judge method for this dataset.

## D Human Annotation Details

We recruit three undergraduate students from UCSB with no prior exposure to this study. The undergraduate students are Asian or Asian American, and all of them are fluent in English. In total, we collected  $6 \times 50 \times 12 \times 3 = 10,800$  human annotations across different datasets and models. Each annotator spends approximately six hours on the task and receives a \$75 Amazon gift card, which is the maximum compensation permitted by the department. The annotators are informed that their work will be included in a research paper, but the topic of the paper is not disclosed in order to minimize potential bias.

Inter-annotator agreement is high (Table 4), indicating the overall quality and consistency of human annotations. In the few cases of disagreement, we apply majority voting to determine the final label.

Furthermore, human annotations (after majority voting) exhibit strong agreement with the LLM-judge (Table 4). This demonstrates that the LLM-judge provides reliable correctness judgments when comparing model answers against ground truth across diverse prompt templates.

The human annotation instructions we use are shown below (from “Hi” to “[GOOGLE SHEETS LINK]”).

Hi [NAME],  
 I have shared a Google spreadsheet with you. Your task would be comparing the model’s answer (“pred” column) with the ground truth (“gt” column), given the question (“question” column). Put 1 in the “human\_evaluation” column if the model’s answer is correct. Put 0 in the “human\_evaluation” column if the model’s answer is wrong. If the

| Dataset       | Human-Human | Human-LLM |
|---------------|-------------|-----------|
| Arc-Challenge | 0.9856      | 0.9785    |
| OpenbookQA    | 0.9919      | 1.0000    |
| GPQA-Diamond  | 0.9083      | 0.9778    |
| NarrativeQA   | 0.6870      | 0.6699    |
| MATH          | 0.8943      | 0.8563    |
| SimpleQA      | 0.7812      | 0.9849    |
| Combined      | 0.9010      | 0.9247    |

Table 4: Human annotation results show strong consistency across annotators. After majority voting, they also align closely with LLM-as-a-Judge evaluations. We report Fleiss’  $\kappa$  (Fleiss, 1971) for agreement among three human annotators (Human–Human) and Cohen’s  $\kappa$  (Cohen, 1960) for agreement between majority-voted human annotations and LLM-as-a-Judge results (Human–LLM).

|                     | Llama-2 | Mistral |
|---------------------|---------|---------|
| Heuristic Mean      | 0.3706  | 0.6125  |
| Heuristic Std       | 0.1413  | 0.0389  |
| LLM-as-a-Judge Mean | 0.5587  | 0.6291  |
| LLM-as-a-Judge Std  | 0.0182  | 0.0098  |

Table 5: Accuracy scores (Mean and Std) across paraphrased prompts for older LLMs.

model’s answer is ambiguous, put 0.  
[GOOGLE SHEETS LINK]

## E Additional Results on Older LLMs

To examine whether prompt sensitivity is a persistent issue across model generations, we further evaluated two earlier instruction-tuned LLMs: **Llama-2-7B-Chat** and **Mistral-7B-Instruct-v0.1**.

From Table 5, we find that using LLM-as-a-Judge significantly reduces variance across paraphrased prompts. These additional results indicate that the reduction of prompt-induced variability is not due to recent models “fixing” the issue, but rather that such variability has consistently been an artifact of evaluation methods rather than a fundamental inconsistency in model behavior.