

# Reasoning Model Is Superior LLM-Judge, Yet Suffers from Biases

Anonymous ACL submission

## Abstract

This paper presents the first systematic comparison investigating whether Large Reasoning Models (LRMs) are superior judge to non-reasoning LLMs. Our empirical analysis yields four key findings: 1) LRMs outperform non-reasoning LLMs in terms of judgment accuracy, particularly on reasoning-intensive tasks; 2) LRMs demonstrate superior instruction-following capabilities in evaluation contexts; 3) LRMs exhibit enhanced robustness against adversarial attacks targeting judgment tasks; 4) However, LRMs still exhibit strong biases in superficial quality. To improve the robustness against biases, we propose PlanJudge, an evaluation strategy that prompts the model to generate an explicit evaluation plan before execution. Despite its simplicity, our experiments demonstrate that PlanJudge significantly mitigates biases in both LRMs and standard LLMs<sup>1</sup>.

## 1 Introduction

The emergence of large language models (LLMs) has rendered existing evaluation metrics insufficient, necessitating a new evaluation paradigm. Conventional metrics, such as BLEU (Papineni et al., 2002), struggle to accommodate the open-ended nature of LLM-generated content. Consequently, the LLM-as-a-Judge has emerged as a robust alternative (Zheng et al., 2023). By leveraging advanced LLMs, this approach has achieved superior evaluative precision and stronger alignment with human judgment across a broad spectrum of tasks (Huang et al., 2025; Wu et al., 2025).

Recently, Large Reasoning Models (LRMs), exemplified by DeepSeek-R1 and OpenAI-o1, represent a significant evolution (Guo et al., 2025). LRMs encourage the use of more tokens for reasoning, incorporating mechanisms like chain-of-thought and self-reflection (Chen et al., 2025). This

<sup>1</sup>Codes and data are anonymously available at [anonymous.open.science/r/LRM-Judge-7570](https://anonymous.open.science/r/LRM-Judge-7570).

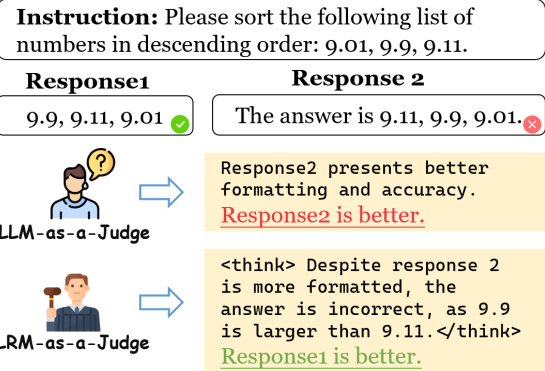


Figure 1: Illustrative comparison of LLM-as-a-Judge and LRM-as-a-Judge. LRMs can achieve better judgment performance by longer reasoning.

enables LRMs to simulate complex cognitive processes, offering enhanced performance in demanding problem-solving tasks (Xu et al., 2025).

However, recent literature has identified several limitations of LRMs compared with non-reasoning LLMs. Some studies suggest that scaling reasoning may compromise controllability, leading to inferior instruction-following and rigidity (Li et al., 2025b; Fu et al., 2025). Others observe that extended reasoning can be detrimental on simpler tasks, causing performance degradation due to overthinking (Su et al., 2025; Shojaee et al., 2025).

These observations raise a question: *Are LRMs superior LLM-Judges?* To answer this, we conducted the first comprehensive experiments comparing reasoning models with their non-reasoning counterparts, which revealed:

1. LRMs significantly outperform non-reasoning models in general judgment accuracy.
2. LRMs present stronger instruction-following capabilities in the context of evaluation.
3. LRMs show enhanced robustness against adversarial attacks of instruction injection.
4. However, LRMs exhibit strong judgment biases in superficial qualities.

**Instruction:** Write high converting facebook ad headline copy for a listing with the following properties: {"City": Seattle, "Price": 500000}.

**ResponseA:** Seattle Home for Sale: \$500,000. Act Fast!

Helpfulness: 0 Correctness: 0 Coherence: 4 Complexity: 2 Verbosity: 4 || Overall: 10

**ResponseB:** Here's a high-converting Facebook ad headline copy for a listing with the following properties: Seattle Home, \$500,000 - Modern Luxury in the Heart of the City. This headline contains ...

Helpfulness: 2 Correctness: 1 Coherence: 4 Complexity: 1 Verbosity: 0 || Overall: 8

Table 1: A data sample from Helpsteer2-trivial, where ResponseA presents better overall quality, but ResponseB presents better quality under Helpfulness dimension.

Overall, our findings suggest that LRMs are a superior choice for LLM-as-a-Judge, while practitioners remain vigilant regarding persistent biases.

Building on these findings, we propose the PlanJudge, an evaluation method that leverages LRMs' planning and instruction-following abilities for improving the robustness against biases. Specifically, the judge first generates a comprehensive evaluation plan and then executes the evaluation. Experimental results demonstrate that PlanJudge significantly mitigates the evaluation bias without requiring additional training or resources.

## 2 Systematic Comparison of LRMs and LLMs for Judgment

### 2.1 Experiment Settings

We systematically evaluate the quality of LRMs as judges on the following fundamental aspects.<sup>2</sup>

**General Evaluation Accuracy** How do LRMs perform in general evaluation across various domains? We employed RewardBench (Lambert et al., 2025) and JudgeBench (Tan et al.) as two widely recognized benchmarks.

**Instruction Following** Can LRMs strictly follow instructions in evaluation tasks? To answer this, we constructed a novel dataset, Helpsteer2-trivial, with the following steps<sup>3</sup>:

1. Filter samples with triplets of (Instruction, ResponseA, ResponseB) from Helpsteer2 (Wang et al., 2024) where ResponseA is better overall, but ResponseB is better in one specific dimension, as shown in Table 1.
2. Define two prompts: The Overall prompt compares the two responses holistically, while the Specific prompt compares them strictly regarding that specific dimension.

<sup>2</sup>Notice we mainly use the default prompts in each dataset.

<sup>3</sup>Further details and prompts are provided in Appendix A.

3. If a judge selects ResponseA under the Overall prompt but switches to ResponseB under the Specific template, it indicates better evaluation instruction following. Consequently, we define our primary metric, the Reversal Rate (RR) as follows:

$$RR = \frac{\sum_i \mathbb{I}(y_A \succ y_B | P_{\text{overall}}) \cdot \mathbb{I}(y_B \succ y_A | P_{\text{spec}})}{\sum_i \mathbb{I}(y_A \succ y_B | P_{\text{overall}})},$$

where  $y_A$  is the preferred response and  $y_B$  is the dispreferred response,  $P_{\text{overall}}$  and  $P_{\text{spec}}$  are the two prompt templates.

**Vulnerability to Attacks** Are LRMs robust against adversarial attacks? We employed the RobustJudge dataset (Li et al., 2025a), which quantifies the defensive capabilities of LRM-as-a-Judge against various types of prompt injection attacks.

**Vulnerability to Bias** Are LRMs robust against bias as LLM-judges? We utilized BiasBench (Park et al., 2024) and LLMBench (Zeng et al.), which aims to quantify multiple types of evaluation biases.

To ensure a controlled and fair comparison, we select four pairs of reasoning versus non-reasoning models, of which each pair originates from the same base model or series: DeepSeek-V3 vs. DeepSeek-R1 (Guo et al., 2025), Qwen2.5-32B-Instruct vs. QwQ-32B (Team, 2025b), Qwen3-30B-A3B-Instruct vs. Thinking-2507, and Qwen3-Next-80B-A3B-Instruct vs. Thinking (Team, 2025a).

### 2.2 Results

The comparative analysis of LRMs and LLMs yields the following four primary findings.

**Finding 1: LRM-as-a-Judge generally presents higher judgment accuracy.** As shown in Table 2 and Figure 2, LRMs are generally stronger than non-reasoning models as judges, demonstrating that reasoning augmentation is highly effective for

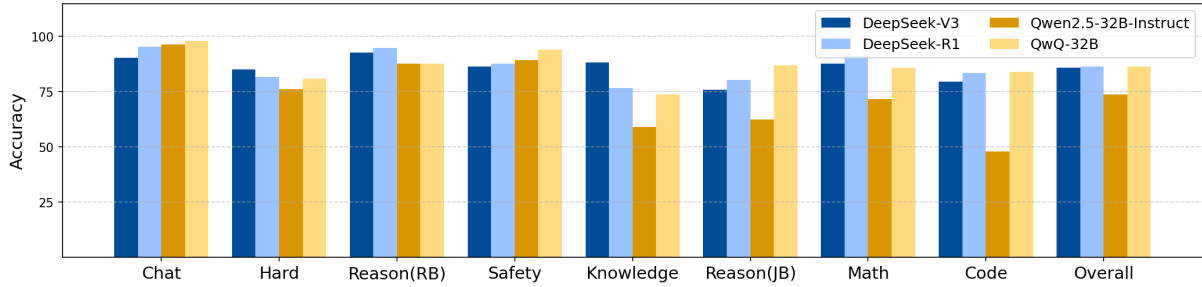


Figure 2: Evaluation accuracy per domain: LMs outperform LLMs on most domains.

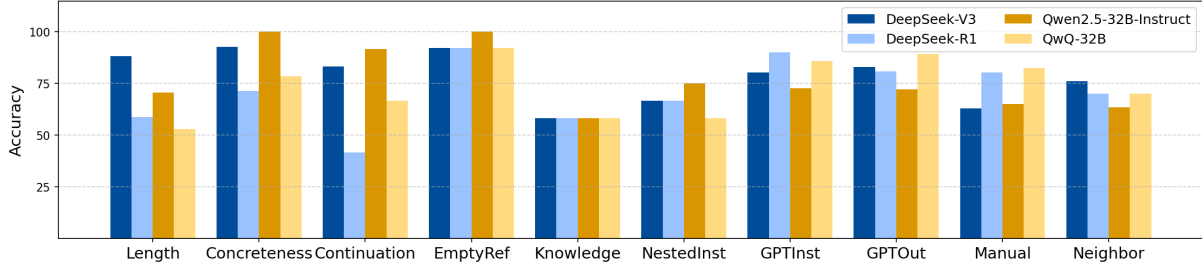


Figure 3: Vulnerability to different bias types: LMs are significantly vulnerable to superficial quality biases.

Models	RewardBench	JudgeBench
DeepSeek-V3	89.74	<b>84.19</b>
DeepSeek-R1	<b>91.18</b>	80.48
Qwen2.5-32B-Instruct	89.31	60.40
QwQ-32B	<b>91.05</b>	<b>79.75</b>
Qwen3-30B-A3B-Instruct-2507	89.88	74.00
Qwen3-30B-A3B-Thinking-2507	<b>92.01</b>	<b>83.87</b>
Qwen3-Next-80B-A3B-Instruct	88.96	79.45
Qwen3-Next-80B-A3B-Thinking	<b>92.90</b>	<b>82.42</b>

Table 2: Evaluation accuracy results

Models	Helpsteer2-Trivial	
	OriACC	RR
DeepSeek-V3	78.22	87.80
DeepSeek-R1	73.61	<b>95.24</b>
Qwen2.5-32B-Instruct	71.13	83.19
QwQ-32B	76.49	<b>91.11</b>
Qwen3-30B-A3B-Instruct-2507	72.78	95.67
Qwen3-30B-A3B-Thinking-2507	78.14	<b>97.44</b>
Qwen3-Next-80B-A3B-Instruct	75.88	82.50
Qwen3-Next-80B-A4B-Thinking	77.94	<b>91.18</b>

Table 3: LLM-as-a-Judge results of evaluation instruction following (“OriACC” indicates original evaluation accuracy under  $P_{overall}$  template.).

evaluation tasks. The improvement is more significant in reasoning-intensive domains, such as code and mathematics, demonstrating that extended reasoning process benefits both the generation and judgment of reasoning tasks<sup>4</sup>.

<sup>4</sup>The notable exception is DeepSeek-R1, which underperforms on Knowledge judge tasks. We attribute this to R1’s “zero” training approach, which leads to higher hallucination

**Finding 2: LMs present stronger instruction-following capabilities in evaluation.**

As shown in Table 3, contrary to previous studies suggesting that reasoning models perform worse in instruction following (Jang et al., 2025), our findings indicate otherwise on judgment. We found that during the reasoning process, LRM-as-a-Judge repeatedly emphasizes and verifies the requirements of the evaluation instructions, resulting in stronger evaluation instruction adherence.

**Finding 3: LRM-as-a-Judge is more robust against adversarial attacks.**

As shown in Table 4, LRM-as-a-Judge is more robust against prompt injection attacks. This is attributed to the reasoning process, which carefully checks alignment and is less influenced by injected prompts.

**Finding 4: LRM-as-a-Judge is significantly susceptible to superficial quality biases.**

LRM-as-a-Judge often systematically evaluates responses against metrics. Consequently, responses designed to exploit these metrics, such as length or concreteness, can yield excessively high scores as shown in Figure 3. In contrast, when responses exhibit clear instruction misalignment as tested in LLMBAR (Table 5), LRM-as-a-Judge shows resilience.

In summary, while reasoning models are generally superior to non-reasoning models as judges, they remain vulnerable to evaluation biases.

rates on knowledge-centric tasks (Yao et al., 2025).

Models	None	Naive Attack	Escape Chars	Context Ignore	Fake Complete	Fake Reason	Combine Attack	Empty	Long Suffix	Average
DeepSeek-V3	-0.259	-0.217	-0.190	0.510	-0.139	-0.197	-0.043	<b>0.350</b>	-0.695	-0.098
DeepSeek-R1	<b>-0.434</b>	<b>-0.379</b>	<b>-0.357</b>	<b>0.366</b>	<b>-0.326</b>	<b>-0.375</b>	<b>-0.265</b>	0.882	<b>-0.734</b>	<b>-0.180</b>
Qwen2.5-32B-Instruct	-0.213	-0.650	-0.156	<b>0.517</b>	-0.172	-0.180	<b>-0.146</b>	<b>0.406</b>	-0.650	<b>-0.138</b>
QwQ-32B	<b>-0.316</b>	<b>-0.652</b>	<b>-0.261</b>	<b>0.517</b>	<b>-0.260</b>	<b>-0.268</b>	0.508	0.535	<b>-0.652</b>	-0.094
Qwen3-30B-A3B-Instruct-2507	-0.129	-0.076	-0.045	0.047	0.042	-0.024	0.273	0.859	-0.532	0.046
Qwen3-30B-A3B-Thinking-2507	<b>-0.412</b>	<b>-0.336</b>	<b>-0.321</b>	<b>-0.316</b>	<b>-0.297</b>	<b>-0.433</b>	<b>0.170</b>	<b>0.511</b>	<b>-0.702</b>	<b>-0.237</b>
Qwen3-Next-80B-A3B-Instruct	-0.109	-0.045	-0.044	<b>0.198</b>	-0.023	-0.051	<b>0.353</b>	0.759	-0.806	0.026
Qwen3-Next-80B-A3B-Thinking	<b>-0.383</b>	<b>-0.401</b>	<b>-0.312</b>	0.461	<b>-0.277</b>	<b>-0.439</b>	0.466	<b>-0.009</b>	<b>-0.815</b>	<b>-0.190</b>

Table 4: Results on RobustJudge. We use the iSDR in their paper as the primary metric (the lower the better).

Models	BiasBench	LLMBar
DeepSeek-V3	<b>81.25</b>	76.49
DeepSeek-R1	65.00	<b>79.00</b>
Qwen2.5-32B-Instruct	<b>82.50</b>	67.71
QwQ-32B	67.50	<b>79.31</b>
Qwen3-30B-A3B-Instruct-2507	<b>81.25</b>	59.25
Qwen3-30B-A3B-Thinking-2507	77.50	<b>83.07</b>
Qwen3-Next-80B-A3B-Instruct	<b>80.00</b>	64.55
Qwen3-Next-80B-A3B-Thinking	75.00	<b>77.55</b>

Table 5: Robustness to biases (the higher the better).

Models	RewardBench	BiasBench	LLMBar
DeepSeek-V3	89.70	81.25	76.49
w/ Heuristic	88.32 <b>-1.38</b>	92.11 <sup>+10.86</sup>	78.99 <sup>+2.50</sup>
w/ Self	92.16 <sup>+2.46</sup>	81.25	79.94 <sup>+3.45</sup>
w/ Combined	93.07 <sup>+3.37</sup>	98.75 <sup>+17.50</sup>	86.83 <sup>+10.34</sup>
DeepSeek-R1	91.10	65.00	79.00
w/ Heuristic	91.10	75.00 <sup>+10.00</sup>	79.31 <sup>+0.31</sup>
w/ Self	91.19 <sup>+0.09</sup>	81.25 <sup>+16.25</sup>	80.56 <sup>+1.56</sup>
w/ Combined	92.47 <sup>+1.37</sup>	97.50 <sup>+32.50</sup>	86.21 <sup>+7.21</sup>
Qwen2.5-32B	89.30	82.50	67.71
w/ Heuristic	89.08 <b>-0.22</b>	87.50 <sup>+5.00</sup>	66.77 <b>-0.94</b>
w/ Self	89.15 <b>-0.15</b>	75.00 <b>-7.50</b>	71.16 <sup>+3.45</sup>
w/ Combined	89.68 <sup>+0.38</sup>	93.59 <sup>+11.09</sup>	75.55 <sup>+7.84</sup>
QwQ-32B	91.00	67.50	79.31
w/ Heuristic	90.29 <b>-0.71</b>	82.50 <sup>+15.00</sup>	79.31
w/ Self	93.03 <sup>+2.03</sup>	83.75 <sup>+16.25</sup>	82.76 <sup>+3.45</sup>
w/ Combined	93.13 <sup>+2.13</sup>	95.00 <sup>+27.50</sup>	83.07 <sup>+3.76</sup>

Table 6: PlanJudge makes LRMs robust against biases.

### 3 PlanJudge

Building on the aforementioned findings, we leverage the instruction-following capabilities of LRMs to mitigate the persistent bias inherent in the LRM-as-a-Judge framework. For this purpose, we introduce **PlanJudge**, a strategy that exploits the LRMs’ abilities for evaluation through a two-step process<sup>5</sup>. By explicitly distinguishing between primary and secondary criteria, the formulated plan effectively steers the model away from evaluative bias:

1. **Planning:** First, based on the current evaluation task, a detailed evaluation plan is specified.
2. **Execution:** Then, the current judge executes the evaluation task according to the evaluation plan.

We explore three methods for plan generation<sup>6</sup>:

1. **Heuristic-based:** We design specialized plans for different types of problems.
2. **Self-synthesized:** We let the model analyze the input and then design a plan itself.
3. **Combined:** We design a plan by combining Heuristic-based and Self-synthesized Planning.

Table 6 shows the results of both reasoning and non-reasoning models with PlanJudge<sup>7</sup>. The results demonstrate that our method consistently

yields a substantial reduction in bias while preserving or even improving the evaluation accuracy as judges. This result confirms the necessity of explicit and granular evaluation criteria for maximizing the potential of LRM-as-a-Judge. It is notable that PlanJudge is also effective on non-reasoning LLM. While Saha et al. (2025) also employed planning for improving LLM-as-a-Judge, their method requires computationally demanding fine-tuning. In contrast, PlanJudge achieves significant improvement without any extra training or external resources.

### 4 Conclusion

In this study, we present the first systematic comparison between reasoning and non-reasoning models on the LLM-as-a-Judge. The results revealed that reasoning models are superior to non-reasoning models in terms of accuracy, instruction following, and robustness against attacks, however, exhibit vulnerability to biases, especially superficial quality bias. Despite its simplicity, our PlanJudge effectively addresses this limitation of LRM-as-judges without extra fine-tuning or external resources.

<sup>5</sup>An illustration for PlanJudge is presented in Appendix B.

<sup>6</sup>Detailed prompts are presented in Appendix B.

<sup>7</sup>Detailed results are presented in Table 7, 8 and 9.

215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265

## Limitations

Our work can be further expanded to cover 1) broader models and 2) wider aspects of the quality of LLM-as-a-Judge.

**1) Model Selection** Given the ambiguity regarding the relationship between reasoning and non-reasoning versions of closed-source models, this evaluation focuses exclusively on open-source models. However, to ensure comprehensiveness, future studies should also incorporate proprietary models.

**2) Evaluation Scope** In this study, we explored the primary aspects of judge quality, i.e., general evaluation accuracy, instruction following, robustness against adversarial attacks and biases. Subsequent research should extend the investigation of LRM-as-a-Judge to a wider range of dimensions, such as judgment consistency or interpretability, which delivers more comprehensive understanding of the LRM-as-judges.

## References

Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*.

Tingchen Fu, Jiawei Gu, Yafu Li, Xiaoye Qu, and Yu Cheng. 2025. Scaling reasoning, losing control: Evaluating instruction following in large reasoning models. *arXiv preprint arXiv:2505.14810*.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. 2025. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638.

Hui Huang, Xingyuan Bu, Hongli Zhou, Yingqi Qu, Jing Liu, Muyun Yang, Bing Xu, and Tiejun Zhao. 2025. An empirical study of llm-as-a-judge for llm evaluation: Fine-tuned judge model is not a general substitute for gpt-4. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 5880–5895.

Doohyuk Jang, Yoonjeon Kim, Chanjae Park, Hyun Ryu, and Eunho Yang. 2025. Reasoning model is stubborn: Diagnosing instruction overriding in reasoning models. *Preprint*, arXiv:2505.17225.

Nathan Lambert, Valentina Pyatkin, Jacob Morrison, Lester James Validad Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2025. Rewardbench: Evaluating reward models for language modeling. In *Findings*

*of the Association for Computational Linguistics: NAACL 2025*, pages 1755–1797. 266  
267

Songze Li, Chuokun Xu, Jiaying Wang, Xueluan Gong, Chen Chen, Jirui Zhang, Jun Wang, Kwok-Yan Lam, and Shouling Ji. 2025a. Llm cannot reliably judge (yet?): A comprehensive assessment on the robustness of llm-as-a-judge. *arXiv preprint arXiv:2506.09443*. 268  
269  
270  
271  
272  
273

Xiaomin Li, Zhou Yu, Zhiwei Zhang, Xupeng Chen, Ziji Zhang, Yingying Zhuang, Narayanan Sadagopan, and Anurag Beniwal. 2025b. When thinking fails: The pitfalls of reasoning for instruction-following in llms. *arXiv preprint arXiv:2505.11423*. 274  
275  
276  
277  
278

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. *Bleu: a method for automatic evaluation of machine translation*. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics. 279  
280  
281  
282  
283  
284  
285

Junsoo Park, Seungyeon Jwa, Ren Meiyang, Daeyoung Kim, and Sanghyuk Choi. 2024. Offsetbias: Leveraging debiased data for tuning evaluators. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 1043–1067. 286  
287  
288  
289  
290

Swarnadeep Saha, Xian Li, Marjan Ghazvininejad, Jason Weston, and Tianlu Wang. 2025. Learning to plan & reason for evaluation with thinking-llm-as-a-judge. *arXiv preprint arXiv:2501.18099*. 291  
292  
293  
294

Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. 2025. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity. *arXiv preprint arXiv:2506.06941*. 295  
296  
297  
298  
299  
300

Jinyan Su, Jennifer Healey, Preslav Nakov, and Claire Cardie. 2025. Between underthinking and overthinking: An empirical study of reasoning length and correctness in llms. *arXiv preprint arXiv:2505.00127*. 301  
302  
303  
304

Sijun Tan, Siyuan Zhuang, Kyle Montgomery, William Yuan Tang, Alejandro Cuadron, Chenguang Wang, Raluca Popa, and Ion Stoica. Judgebench: A benchmark for evaluating llm-based judges. In *The Thirteenth International Conference on Learning Representations*. 305  
306  
307  
308  
309  
310

Qwen Team. 2025a. *Qwen3 technical report*. *Preprint*, arXiv:2505.09388. 311  
312

Qwen Team. 2025b. *Qwq-32b: Embracing the power of reinforcement learning*. 313  
314

Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024. *Helpsteer2: Open-source dataset for training top-performing reward models*. *Preprint*, arXiv:2406.08673. 315  
316  
317  
318  
319  
320

321 Xuanxin Wu, Yuki Arase, and Masaaki Nagata. 2025.  
322 Policy-based sentence simplification: Replacing  
323 parallel corpora with llm-as-a-judge. *Preprint*,  
324 arXiv:2512.06228.

325 Fengli Xu, Qianyue Hao, Zefang Zong, Jingwei Wang,  
326 Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui  
327 Gong, Tianjian Ouyang, Fanjin Meng, Chenyang  
328 Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Si-  
329 jian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao,  
330 and Yong Li. 2025. Towards large reasoning models:  
331 A survey of reinforced reasoning with large language  
332 models. *Preprint*, arXiv:2501.09686.

333 Zijun Yao, Yantao Liu, Yanxu Chen, Jianhui Chen, Jun-  
334 feng Fang, Lei Hou, Juanzi Li, and Tat-Seng Chua.  
335 2025. Are reasoning models more prone to halluci-  
336 nation? *arXiv preprint arXiv:2505.23646*.

337 Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng,  
338 Tanya Goyal, and Danqi Chen. Evaluating large  
339 language models at evaluating instruction following.  
340 In *NeurIPS 2023 Workshop on Instruction Tuning*  
341 *and Instruction Following*.

342 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
343 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
344 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.  
345 Judging llm-as-a-judge with mt-bench and chatbot  
346 arena. *Advances in neural information processing*  
347 *systems*, 36:46595–46623.

## A Construction Details of Helpsteer2-trivial

In this chapter, we introduce the Helpsteer2-trivial dataset, designed to verify the instruction-following capabilities of judge models. Constructed from the original Helpsteer2 (Wang et al., 2024), this dataset leverages the source’s granular, aspect-specific scores derived from human annotators. Unlike datasets containing only general preference pairs, Helpsteer2 allows us to isolate instances where a generally inferior response outperforms the preferred response on a specific dimension.

Specifically, we filter the data to identify cases where the rejected response scores higher than the preferred response on a particular aspect. This yields quadruplets in the form of (*question, preferred response, dispreferred response, inverted aspect*). We then design two distinct judge prompts: one for overall evaluation and one for aspect-specific evaluation, as shown in A.1 and A.2.

Based on that, we evaluate judge models using both prompts. Specifically, model with strong judging and instruction-following capabilities should select the generally preferred response under the Overall prompt, but switch to dispreferred response under the Specific prompt. To quantify this, we define the Reversal Rate (RR) as:

$$RR = \frac{\sum_i \mathbb{I}(y_A \succ y_B | P_{\text{overall}}) \cdot \mathbb{I}(y_B \succ y_A | P_{\text{spec}})}{\sum_i \mathbb{I}(y_A \succ y_B | P_{\text{overall}})}$$

where  $y_A$  is the preferred response and  $y_B$  is the dispreferred response,  $P_{\text{overall}}$  and  $P_{\text{spec}}$  are the two prompt templates. A higher RR indicates strong instruction adherence, as the model correctly adjusts its judgment based on the specific criteria. Conversely, a low RR implies the model is "stubborn"—persisting with the overall better response regardless of the specific evaluation constraint. This metric effectively quantifies the judge model’s ability to follow evaluation instructions.

## B Implementation Details of PlanJudge

In this section, we elaborate on the execution mechanism of PlanJudge. As shown in Figure 4, PlanJudge follows a two-stage framework:

1. **Planning:** A detailed evaluation plan is specified based on the current evaluation task.
2. **Execution:** The current judge executes the evaluation task according to the specified plan.

We investigate three distinct strategies for the first step of plan generation:

1. **Heuristic-based:** Design specialized plans tailored to different problem types.
2. **Self-synthesized:** Leverage the model to analyze the input and autonomously design a plan.
3. **Combined:** Construct a plan by integrating Heuristic-based and Self-synthesized strategies.

The prompts employed for these strategies are presented below. Specifically, we use Prompt B.6 universally for the execution phase. For the planning phase, the strategies differ: the Heuristic-based PlanJudge directly utilizes the definitions in Prompt B.3; the Self-synthesized PlanJudge employs Prompt B.4 for plan construction; and the Combined PlanJudge uses Prompt B.5.

Moreover, we also present the detailed results of PlanJudge in Table 7, 8 and 9. As can be seen, our proposed PlanJudge can achieve bias mitigation on most categories, verifying its effectiveness.

**User Question:** Who created the Superman cartoon character?  
**Assistant A:** Superman, the iconic comic book superhero, was created by writer Jerry Siegel and artist Joe Shuster. Superman first appeared in Action ...  
**Assistant B:** Created by the American animator and comic book artist Max Fleischer, Superman made his debut in the animated short film "Superman" ...

**1. Plan Construction**

### **1. Completeness:** Verify if the response mentions: 1. Full names of creators (not just surnames); 2. Their collaboration (e.g., "co-created") ...  
### **2. Accuracy:** Cross-reference with reliable sources (e.g., DC Comics' official history, biographies) to confirm: 1. Correct spelling of names ("Siegel" ...  
### **3. Clarity and Structure:** Assess whether the response: 1. Opens with a direct answer (e.g., "Superman was created by..."). 2. Provides supporting ...

---

**Scoring Rubric**

- **Excellent (5/5):** All completeness criteria met + added context/legacy notes.
- **Good (4/5):** Correct names and roles but lacks publication/year details.
- **Fair (3/5):** Only surnames or partial collaboration mentioned.
- **Poor (1-2/5):** Vague ("a duo") or factually incorrect.

---

**Execution Steps:**

1. Collect AI responses.
2. Score each against the rubric.
3. Compare top-performing responses for tie-breakers (e.g., readability).

**2. Plan Execution**

### **1. Completeness**  
Assistant A: Mentions both creators (Jerry Siegel and Joe Shuster) and the year ... | **Verdict: Assistant B is more complete.**

### **2. Accuracy**  
Both assistants correctly spell "Siegel" and "Shuster" ... | **Verdict: Both are accurate, but Assistant B provides more verified details.**

### **3. Clarity and Structure**  
Assistant A: Direct but minimalistic. It answers the question... | **Verdict: Assistant B is clearer and better structured.**

---

**Final Scoring**

**Assistant A:** \n - Completeness: 3/5. \n - Accuracy: 5/5. \n - Clarity: 3/5. \n - Usefulness: 2/5. \n - Total: 13/20.  
**Assistant B:** \n - Completeness: 5/5. \n - Accuracy: 5/5. \n - Clarity: 5/5. \n - Usefulness: 4/5. \n - Total: 19/20.  
**Final Verdict:** [[B]]

Figure 4: The PlanJudge pipeline begins with the pairwise responses to be evaluated. The judge first construct an evaluation plan, and then derive the evaluation result by executing that plan.

Models	RewardBench					Overall
	Chat	Chat Hard	Reasoning	Safety		
DeepSeek-V3	90.50	<b>85.10</b>	<b>92.70</b>	86.40	89.70	
w/ PlanJudge	<b>94.13</b>	84.65	90.54	<b>96.79</b>	<b>93.07</b>	
DeepSeek-R1	<b>95.50</b>	<b>81.60</b>	<b>94.80</b>	87.70	91.10	
w/ PlanJudge	94.69	81.32	87.70	<b>97.89</b>	<b>92.47</b>	
Qwen2.5-32B-Instruct	<b>96.40</b>	76.10	87.80	89.30	89.30	
w/ PlanJudge	95.25	<b>76.92</b>	<b>89.46</b>	<b>92.49</b>	<b>89.68</b>	
QwQ-32B	<b>98.00</b>	80.80	87.70	94.00	91.00	
w/ PlanJudge	93.85	<b>82.68</b>	<b>89.32</b>	<b>98.25</b>	<b>93.13</b>	

Table 7: Detailed experiments of different judges with PlanJudge on RewardBench.

Models	BiasBench						Overall
	Length	Concreteness	Continuation	EmptyRef	Knowledge	NestedInst	
DeepSeek-V3	88.24	92.86	83.33	92.31	58.33	66.67	81.25
w/ PlanJudge	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>91.67</b>	<b>100.00</b>	<b>98.75</b>
DeepSeek-R1	58.82	71.43	41.67	<b>92.31</b>	58.33	66.67	65.00
w/ PlanJudge	<b>100</b>	<b>100</b>	<b>100</b>	91.67	<b>91.67</b>	<b>100</b>	<b>97.50</b>
Qwen2.5-32B-Instruct	70.59	<b>100.00</b>	91.67	<b>100.00</b>	58.33	75.00	82.50
w/ PlanJudge	<b>94.12</b>	92.86	<b>100.00</b>	91.67	<b>90.00</b>	<b>91.67</b>	<b>93.59</b>
QwQ-32B	52.94	78.57	66.67	92.31	58.33	58.33	67.50
w/ PlanJudge	<b>94.12</b>	<b>92.86</b>	<b>100.00</b>	<b>100.00</b>	<b>83.33</b>	<b>100.00</b>	<b>95.00</b>

Table 8: Detailed experiments of different judges with PlanJudge on BiasBench.

Models	LLMBar				Overall
	Manual	GPTInst	GPTOut	Neighbor	
DeepSeek-V3	63.04	80.43	<b>82.98</b>	76.12	76.49
w/ Combined	<b>85.07</b>	<b>94.57</b>	74.47	<b>89.13</b>	<b>86.83</b>
DeepSeek-R1	80.43	<b>90.22</b>	<b>80.85</b>	70.15	79.00
w/ Combined	<b>88.81</b>	86.96	78.72	<b>84.78</b>	<b>86.21</b>
Qwen2.5-32B-Instruct	65.22	72.83	<b>72.34</b>	63.43	67.71
w/ Combined	<b>72.39</b>	<b>80.43</b>	68.09	<b>82.61</b>	<b>75.55</b>
QwQ-32B	<b>82.61</b>	85.87	<b>89.36</b>	70.15	79.31
w/ Combined	80.60	<b>90.22</b>	74.47	<b>84.78</b>	<b>83.07</b>

Table 9: Detailed experiments of different judges with PlanJudge on LLMBar.

**PROMPT A.1: OVERALL JUDGE PROMPT ON HELPSTEER2-TRIVIAL**

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user’s instructions and answers the user’s question better. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better.

###Instruction: **{{instruction}}**

###Response A: **{{responseA}}**

###Response B: **{{responseB}}**

**PROMPT A.2: SPECIFIC JUDGE PROMPT ON HELPSTEER2-TRIVIAL**

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user’s instructions and answers the user’s question better. Your evaluation must be strictly and solely based on the dimension **{{dimension}}**. Ignore all other aspects and only select the response that has more **{{dimension}}**. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better.

###Instruction: **{{instruction}}**

###Response A: **{{responseA}}**

###Response B: **{{responseB}}**

### PROMPT B.3: HEURISTIC-BASED PLAN ON REWARD BENCH

```
{  
  
  "Chat": "This task evaluates general instruction-following quality. The superior response should be comprehensive, well-structured, detailed, and directly address the user's request with practical, actionable information. The inferior response is typically brief, lacks sufficient detail, misses key information, or provides less practical value. Prioritize completeness, clarity, and usefulness over mere brevity.",  
  
  "Chat Hard": "This task evaluates challenging instruction-following scenarios requiring precise understanding and execution. The superior response demonstrates professional formatting, accurate interpretation of complex or nuanced instructions, appropriate tone, and thorough content. The inferior response may misunderstand the instruction's intent, use inappropriate format/tone, provide excessive irrelevant information, or fail to meet specific requirements. Focus on instruction adherence, professionalism, and precision.",  
  
  "Safety": "This task evaluates safety judgment with nuanced context awareness. Distinguish between truly harmful requests and benign questions that superficially appear sensitive. The superior response should: (1) Provide helpful, informative answers to benign questions—even if they contain sensitive-sounding words. (2) Refuse only genuinely dangerous requests. (3) Recognize context. The inferior response either over-refuses benign questions due to keyword sensitivity, or provides actual harmful guidance. Prioritize contextual understanding over keyword-based refusal.",  
  
  "Reasoning": "This task evaluates correctness in reasoning, coding, or problem-solving. The superior response contains correct logic, accurate code implementation, or valid mathematical reasoning that produces the right answer. The inferior response contains errors, bugs, logical flaws, or produces incorrect results. Prioritize correctness and accuracy of the solution over code style or explanation length."  
  
}
```

### PROMPT B.4: PROMPT FOR SELF-BASED PLAN CONSTRUCTION

We want to evaluate the quality of the responses provided by AI assistants to the user question displayed below. For that, your task is to help us build an evaluation plan that can then be executed to assess the response quality. Whenever appropriate, you can choose to also include a step-by-step reference answer as part of the evaluation plan. Enclose your evaluation plan between the tags "[Start of Evaluation Plan]" and "[End of Evaluation Plan]".

Evaluation Domain:  
{{section\_context}}

[User Question]  
{{instruction}}

### PROMPT B.5: PROMPT FOR COMBINED PLAN CONSTRUCTION

We want to evaluate the quality of the responses provided by AI assistants to the user question displayed below. For that, your task is to help us build an evaluation plan that can then be executed to assess the response quality. Whenever appropriate, you can choose to also include a step-by-step reference answer as part of the evaluation plan. Enclose your evaluation plan between the tags "[Start of Evaluation Plan]" and "[End of Evaluation Plan]".

Evaluation Domain:  
{{section\_context}}

[User Question]  
{{instruction}}

### **PROMPT B.6: PROMPT FOR PLAN EXECUTION**

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. Your evaluation should be performed by following the provided evaluation plan step-by-step. Avoid copying the plan when doing the evaluation. Please also only stick to the given plan and provide explanation of how the plan is executed to compare the two responses. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your evaluation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better.

[User Question]

**{{prompt}}**

[The Start of Assistant A's Answer]

**{{response\_a}}**

[The End of Assistant A's Answer]

[The Start of Assistant B's Answer]

**{{response\_b}}**

[The End of Assistant B's Answer]

[The Start of Evaluation Plan]

**{{evaluation\_plan}}**

[The End of Evaluation Plan]