

Can You Trust LLM Judgments? Reliability of LLM-as-a-Judge

Kayla Schroeder

Department of Statistics

Northwestern University

kaylaschroeder2026@u.northwestern.edu

Zach Wood-Doughty

Department of Computer Science

Northwestern University

zach@northwestern.edu

Abstract

Large Language Models (LLMs) have become increasingly powerful and ubiquitous, but their stochastic nature poses challenges to the reliability of their outputs. While deterministic settings can improve consistency, they do not guarantee reliability, as a single sample from the model’s probability distribution can still be misleading. Building upon the concept of LLM-as-a-judge, we introduce a novel framework for rigorously evaluating the reliability of LLM judgments, leveraging McDonald’s omega. We evaluate the reliability of LLMs when judging the outputs of other LLMs on standard single-turn and multi-turn benchmarks, simultaneously investigating the impact of temperature on reliability. By analyzing these results, we demonstrate the limitations of fixed randomness and the importance of considering multiple samples, which we show has significant implications for downstream applications. Our findings highlight the need for a nuanced understanding of LLM reliability and the potential risks associated with over-reliance on single-shot evaluations. This work provides a crucial step towards building more trustworthy and reliable LLM-based systems and applications.

1 Introduction

In recent years, Large Language Models (LLMs) have experienced rapid advancements and widespread adoption, with applications ranging from marketing to biotechnology to music (Gozalo-Brizuela and Garrido-Merchán, 2023). As reliance upon LLMs grows, ensuring the quality and trustworthiness of their outputs becomes crucial. We often think about LLM quality in terms of accuracy (how often the LLM is correct). However, equally important are reliability, confidence, and calibration; these three distinct concepts are intertwined and essential for building genuine trust in LLM systems.

LLMs can express how certain, or ‘confident’,

they are about their responses. Calibration measures how well this confidence aligns with the LLM’s actual accuracy (Kotelanski et al., 2023; Spiess et al., 2024). A perfectly calibrated LLM that says it’s 90% confident should be correct about 9 out of 10 times. However, even a perfectly calibrated LLM can be unreliable.

Reliability refers to the consistency of an LLM’s outputs. A reliable weather forecasting LLM, for example, should give similar predictions for the same day, even if those predictions aren’t always accurate. An LLM that sometimes predicts sunshine and sometimes a downpour for the same day is unreliable, even if it’s well-calibrated and highly confident. This kind of consistency is what we aim to capture with internal consistency reliability. A reliable LLM can be consistently wrong, but it’s still reliable because its judgments are stable.

Why is reliability so important? Consider a scenario where an LLM is tasked with evaluating the quality of a medical imaging LLM used to diagnose tumors. Even if the LLM is perfectly calibrated (its confidence accurately reflects its diagnostic ability), if it gives different evaluations for the same AI on different runs, even with different, but accurate, confidence scores, doctors won’t know which evaluation to trust. The inconsistency makes the LLM unreliable and ultimately untrustworthy, regardless of its calibration and confidence.

A singular output from an LLM is the result of a draw from a probabilistic distribution (Vaswani, 2017). Parameters that vary the resulting output, namely temperature and top-k, are increased to obtain more novel or “creative” results - a desirable quality of LLMs (Peeperkorn et al., 2024). However, this inherent stochastic nature raises concerns about individual outputs. Given that a single output from an LLM is a representation of only a single draw from the model’s distribution, trustworthiness of an LLM output is of utmost importance (Mittelstadt et al., 2023). Confidence and calibration

are commonly used within NLP to understand this variability (Lin et al., 2023; Tian et al., 2023), but they don't address reliability directly.

Many research works aim to circumvent this issue by setting a fixed seed and using deterministic settings for the temperature and top-k parameters (Ouyang et al., 2023; Wei et al., 2024; Atil et al., 2024). These studies argue that if an LLM consistently produces the same output under these conditions, it can be considered reliable. However, consistent replication does not guarantee the reliability of the generated text.

Even with deterministic settings, a single LLM output remains a sample from the model's probability distribution, subject to inherent randomness. This results in "fixed randomness," which can lead to significant limitations (Hellrich and Hahn, 2016). Consider our prior example of an LLM evaluator of a medical imaging LLM for tumor diagnosis. If the evaluating LLM, due to its own inherent randomness, consistently fails to identify a critical weakness in the target LLM in a single evaluation run, grave repercussions can arise. Eradicating fixed randomness is crucial to ensure model consistency and reduce the risk of these systematic errors. Internal consistency reliability addresses this problem by assessing how consistently the LLM applies its evaluation criteria across multiple evaluations, rather than relying on a single, potentially flawed output.

The ability to trust LLM outputs is of further importance as LLMs are increasingly used as evaluators (Desmond et al., 2024; Chan et al., 2023a). The LLM-as-a-judge paradigm, coined by Zheng et al. (2023), relies upon the internal consistency reliability of LLM judges. Widespread utilization of LLMs as judges has garnered concern. The ACL Rolling Review Reviewer Guidelines highlight this issue, voicing explicit concern about LLM-only evaluation without validation by listing the methodological issue "If LLMs are used as automated evaluators, is their reliability in this context sufficiently validated?"¹. This underscores the critical need for trust in LLM outputs and their reliability. Reliability here, however, is not explicitly defined, and inconsistencies in how reliability of LLMs should be quantified are present throughout the literature.

While some research has focused on determining LLM reliability by comparing outputs to ground

truth labels (Zhou et al., 2024; Fu et al., 2023a), this approach is often limited by the availability of accurate and comprehensive ground truth data. As LLMs are increasingly employed in complex and subjective domains, such as evaluating a medical imaging LLM for tumor diagnosis where even expert human diagnoses can vary, establishing a reliable ground truth becomes increasingly challenging. In current practice, LLM judgment models are often treated as raters, and inter-rater reliability is calculated on their outputs to assess reliability (Kollitsch et al., 2024; Wang et al., 2024). This metric, as we will show, is an insufficient quantification of reliability as the metric does not take into account the additional complexity of the LLM judge being a model itself and containing randomness.

The unique characteristics of LLMs as evaluators necessitate a more nuanced approach to assessing their reliability. We define reliability in this context as internal consistency reliability, a well-established statistical concept and the most common type of reliability employed across domains Henson (2001). While internal consistency is a well-established concept, its application to the LLM-as-a-judge paradigm, particularly with a focus on seed variation, has been largely overlooked. Internal consistency reliability measures how well different parts of an evaluation (individual LLM judgments) measure the same underlying construct (the LLM's "true" judgment). Internal consistency reliability ensures that the LLM applies its evaluation criteria consistently across various inputs and conditions. Metrics like Cronbach's alpha and McDonald's omega provide quantifiable measures of this internal consistency, allowing us to assess the reliability of an LLM's evaluations by examining the consistency of its responses across a range of prompts or tasks (Cronbach, 1951; McDonald, 2013; Malkewitz et al., 2023).

We introduce a novel framework that leverages McDonald's omega to address a critical gap in current LLM evaluation methodologies by rigorously quantifying the stability and consistency of LLM judgments across multiple replicated evaluations of the same responses. Critically, this work pioneers the systematic investigation of the variability of an LLM evaluator's judgments solely by varying the random seed, holding all other parameters constant. We assess how consistently LLMs apply their evaluation criteria across replications (i.e., multiple independent evaluations of the same responses with fixed parameters), investigating the impact of di-

¹As of December 2024 at <https://aclrollingreview.org/reviewerguidelines>

verse question formats and varying difficulty levels using established benchmarks (BBH, SQuAD, MT-Bench). Our findings provide novel insights into the sensitivity of LLM judgments, revealing potential limitations in their reliability and highlighting the importance of considering this variability in high-stakes applications.²

2 Related Work

Reliability is a cornerstone of NLP research, particularly when human judgment is involved. Studies utilizing human labelers have long emphasized inter-annotator agreement (Bhowmick et al., 2008; Nowak and Rüger, 2010; Amidei et al., 2019). Similarly, topic modeling has extensively explored reliability through various methods, including similarity measures, statistical techniques, and domain expertise (Rieger et al., 2024; Schroeder and Wood-Doughty, 2024; Chuang et al., 2015).

Internal consistency reliability itself, especially Cronbach's alpha, has been heavily utilized as an assessor of question quality in LLM contexts (Shang et al., 2024; Chan et al., 2023b; Biri et al., 2023). Bhandari et al. (2024) uses Cronbach's alpha to compare LLM-generated questions with human-authored ones; Shi et al. (2023) uses Cronbach's alpha to assess the reliability of questions designed to evaluate LLMs. However, these studies focus on the reliability of questions, not the reliability of LLM judgments within the LLM-as-a-judge paradigm, a key distinction in our work.

The concept of reliability appears widely in the LLM literature, but is not consistently defined. Some studies equate reliability with human-model agreement, requiring costly human evaluations (Sun et al., 2024; Zhang et al., 2024). Others rely on accuracy metrics, which necessitate ground truth labels (Zheng et al., 2024; Macherla et al., 2023; Wei et al., 2024). As the complexity and subjectivity of judgment tasks increase, obtaining reliable ground truth becomes increasingly challenging. Semantic consistency has also been proposed as a reliability measure (Govil et al., 2024; Raj et al., 2023), but this approach lacks a strong theoretical foundation and has potential pitfalls, as discussed in Schroeder and Wood-Doughty (2024). Many studies also use the term "reliability" without any rigorous definition (Tan et al., 2024; Liao et al., 2023; Gupta et al., 2024). This lack of a clear and

consistent definition highlights the need for a more principled approach to measuring LLM judgment reliability.

A common approach for establishing reliability across LLM judgment models is the employment of inter-rater reliability (Kollitsch et al., 2024; Wang et al., 2024). This metric, as we will show, is an insufficient quantification of reliability as the metric does not take into account the additional complexity of the LLM judge being a model itself and containing randomness. Beyond inter-rater reliability, the LLM-as-a-judge literature predominantly relies upon accuracy and bias metrics to assess reliability (Chen et al., 2024; Wei et al., 2024; Ye et al., 2024). Some researchers advocate for combining responses from different LLM judgment models to improve reliability (Gu et al., 2024; Patel et al., 2024). To quantify reliability of LLM judgments, human judgments are often used as a benchmark, with methods relying upon agreement with human evaluation and correlation with human-annotated error scores (Jung et al., 2024; Fu et al., 2023b; Dong et al., 2024). Unlike these approaches, which often rely on external benchmarks or human input, our work focuses on the internal consistency of LLM judgments, a critical aspect of reliability that has been largely overlooked in the LLM-as-a-judge paradigm. Our work pioneers the investigation of internal consistency reliability in this context, specifically by isolating the impact of random seed variation, addressing a significant gap in the literature.

3 Limitations of Current Reliability Measures

To illustrate the limitations of inter-rater reliability for LLM judgment reliability, we present a preliminary example. While inter-rater reliability, calculated as the percentage of agreement between raters and detailed further in (Hallgren, 2012), is a useful metric in some contexts, it is not well-suited for the unique characteristics of LLMs as judges (specifically, the inherent randomness within LLM evaluators themselves).

To demonstrate this, we performed an analysis using a simplified version of our Reliability Framework, which is described fully in Section 4. We obtained judgments from three LLM evaluators on responses to BIG-Bench Hard (BBH) questions from Suzgun et al. (2022), using a temperature of 0.25. Each LLM provided 100 judgments for each

²The code for the framework and analyses can be accessed at https://github.com/kaylaschroeder/llm_reliability.

	Min	Q1	Q2	Q3	Max
IRR	0.167	0.395	0.433	0.469	1.000

Table 1: Variation of the inter-rater reliability (IRR) across 3 LLM evaluators across replications.

question, varying only the random seed. We then calculated the inter-rater reliability across the three LLM evaluators for each replication of judgments (one set per replication), for a total of 100 inter-rater reliability values.

Table 1 details the drastic variation in inter-rater reliability, ranging from 0.167 to 1.00. This disconcertingly wide range demonstrates that inter-rater reliability is highly sensitive to random seed variation, making it an unreliable indicator of true judgment quality. This instability in inter-rater reliability underscores the need for a more robust and appropriate measure of LLM judgment reliability, further motivating the development of our proposed (internal consistency) reliability framework.

4 Reliability Framework

Our framework evaluates the reliability of the LLM-as-a-judge paradigm across diverse question formats and difficulty levels. Using LLM responses to benchmark questions, judgment LLMs are prompted repeatedly to select the "best" response based on factors like accuracy, utility, and relevance. The LLM is prompted for evaluation 100 times, varying only the replication while holding other factors constant. All judgment results are then assessed to uncover reliability by applying McDonald's omega.

The reliability framework begins by utilizing the widely accepted benchmarks BIG-Bench Hard (BBH) from [Suzgun et al. \(2022\)](#), SQuAD from [Rajpurkar \(2016\)](#), and MT-Bench from [Zheng et al. \(2023\)](#). These benchmarks offer two key advantages. First, they present diverse question formats: BBH and SQuAD questions are single-turn with specific target answers, while MT-Bench questions are multi-turn and open-ended. Second, they provide varying levels of difficulty, as evidenced by the significantly different accuracy rates of LLMs on BBH and SQuAD questions shown in Table 4. This diversity allows us to assess judgments under a range of conditions.

One question from each category within each benchmark is randomly selected, for a total of 55 questions. Next, five responses for each question

are generated using Chain-of-Thought prompting, as illustrated in Figure 1. The variation in prompting by benchmark is a result of the structuring of each benchmark; BBH responses require a single question, SQuAD responses require a single question as well as context surrounding the question, and MT-Bench requires a two-turn dialogue.

Include step-by-step reasoning in answering the following question:
[Question]

(a) Prompt template for generating responses to BBH questions.

Include step-by-step reasoning in answering the following question:
Context: [Context]
Question: [Question]
Answer:

(b) Prompt template for generating responses to SQuAD questions.

Include step-by-step reasoning in answering the following question:
[Question 1]
[Response 1]
[Question 2]

(c) Prompt template for generating responses to the second turn and final of MT-Bench questions. To generate responses from the first turn, the template is identical to the BBH prompt.

Figure 1: Prompt templates to generate responses from each of the varying benchmark types.

The model responses are then combined in a single prompt, as shown in Figure 2. To prevent judgment bias, the five responses to the respective benchmark question are included in the prompt in a random order as response options [A] - [E]. This random order is held constant across judgment replications, but shuffled for each new benchmark question judgment prompt. In addition to the five responses, the judgment prompt includes the original benchmark question (and context or previous turns if relevant) and prompts the LLM for a judgment of best response including Chain-of-Thought. Here, Chain-of-Thought from [Wei et al. \(2022\)](#) is again utilized as it has been shown to dramatically improve LLM performance [Feng et al. \(2024\)](#); [Suzgun et al. \(2022\)](#); [Chu et al. \(2023\)](#). For the multi-turn responses from MT-Bench, options

[A]-[E] are explicitly labeled as responses to the final turn, as shown in Figure 2b.

Each LLM utilized for judgments is prompted 100 times, varying only the replication, to evaluate which of the responses to a given benchmark question is best (as defined in the prompt) out of the five responses provided. Given the employment of Chain-of-Thought in the initial responses to a given benchmark question, a single best response should exist in response options [A] - [E] for a judge to uncover because each of the five responses contains reasoning in addition to the question response. Similar logic is employed in [Zheng et al. \(2023\)](#) in assessment of judgments as well.

The reliability of the judgments is then calculated using internal consistency reliability, a cornerstone of psychometric analysis and has been extensively studied. Cronbach's alpha, introduced by [Cronbach \(1951\)](#), has been widely used to assess internal consistency, but relies on several assumptions that may not always hold ([Malkewitz et al., 2023; Agbo, 2010](#)). As proposed by [McDonald \(2013\)](#), McDonald's omega addresses these limitations by accounting for a broader range of factor structures. Ultimately, omega provides a more robust measure of internal consistency and has been shown to be more robust to deviations from assumptions than alpha ([Stensen and Lydersen, 2022](#)). In this work, we adopt McDonald's omega as our primary metric, recognizing its superior quantification of reliability.

4.1 McDonald's Omega

McDonald's omega (ω) is formulated as

$$\frac{(\sum_{i=1}^n \lambda_i)^2}{(\sum_{i=1}^n \lambda_i)^2 + \sum_{i=1}^n \theta_{ii} + 2 \sum_{1 \leq i < j \leq n} \theta_{ij}}$$

where n is the number of replications which, in this setting, is the number of times the LLM makes a judgment. The parameter λ_i is the factor loading of the i -th replication. This represents the strength of the relationship between the i -th judgment and the underlying "true" judgment. In simpler terms, factor loadings indicate how much each replication contributes to our understanding of the overall judgment being measured. A higher factor loading suggests that the replication is a strong indicator of the underlying true judgment. The error variance of the i -th replication is given by θ_{ii} , which represents the variability or "noise" associated with the

You are a fair and objective judge tasked with selecting the strongest of the following responses to the provided question. Base your judgment upon the accuracy, utility, and relevance of each. Do not consider length of response, positioning of response or title of response in your judgment. Output the letter of the best response followed by an explanation, and strictly follow the following format: "Best Response: [[letter]]".

Question: [Question]

Responses:

- [A]: [Response A]
- [B]: [Response B]
- [C]: [Response C]
- [D]: [Response D]
- [E]: [Response E]

(a) Prompt template for generating judgments of responses to BBH and SQuAD questions.

You are a fair and objective judge tasked with selecting the strongest of the following responses to the second provided question ("Question 2"). Question 1 is only provided for context. Base your judgment upon the accuracy, utility, and relevance of each. Do not consider length of response, positioning of response or title of response in your judgment. Output the letter of the best response followed by an explanation, and strictly follow the following format: "Best Response: [[letter]]".

Question 1: [Question 1]

Question 2: [Question 2]

Responses:

- [A]: [Response A]
- [B]: [Response B]
- [C]: [Response C]
- [D]: [Response D]
- [E]: [Response E]

(b) Prompt template for generating judgments of responses to MT-Bench questions.

Figure 2: Prompt templates for generating LLM evaluations of responses for each of the varying benchmarks.

i-th judgment. Finally, the covariance between the errors of replications i and j is given by θ_{ij} . This measures the extent to which the errors in different judgments are related. Essentially, McDonald's omega assesses how much of the observed variation in judgments reflects the true underlying value and how much is due to random error or inconsistencies in the measurement process.

To apply omega, we assume:

1. Additive Measurement Error: Error in each LLM judgment replication is independent of the true judgment value. In simpler terms, the error does not systematically over- or underestimate the true judgment. This is a common psychometrics assumption, particularly when dealing with judgments that are expected to be relatively consistent (Alagumalai and Curtis, 2005; Wadkar et al., 2016; Liu et al., 2010).
2. Uncorrelated Errors: Errors in different LLM judgment replications are independent of each other. In other words, the error in one judgment does not influence the error in another judgment. This assumption is trivial in our setting, as each LLM judgment is generated independently so errors are guaranteed to be uncorrelated.
3. Single Latent Trait: All LLM judgments are attempting to measure the same underlying "true" judgment. In our case, each LLM judgment replication aims to identify the most accurate, relevant, and useful response to the same prompt, aligning with this assumption.

5 Experiments and Data

Responses to benchmark questions, using the templates outlined in Figure 1, are generated from five highly popular open-source LLMs: LLaMA-3-8B, Vicuna-7B-v1.5, Gemma-7B, Phi-2, and Falcon-7b (Dubey et al., 2024; Chiang et al., 2023; Team et al., 2024a; Javaheripi et al., 2023; Almazrouei et al., 2023). These models were chosen due to their popularity, with over 100,000 downloads in the prior month on Hugging Face.³ The generated responses are obtained using top-k of 50 and a temperature of 0.75.

Selecting models that are high performers in the Chatbot Arena (Chiang et al., 2024), Starling-LM-

7B-beta, Gemma-1.1-7b-it, and Meta-Llama-3-8B-Instruct are employed for the judgment task Zhu et al. (2023); Team et al. (2024b); AI@Meta (2024). These selected models are all below the threshold of 10B parameters and are top performing models, ranked in the top 100, with rankings of 65 (Meta-Llama-3-8B-Instruct), 78 (Starling-LM-7B-beta), and 98 (Gemma-1.1-7b-it). In addition, we note that the selected models differ from those used for question answering in the prior step. This distinction is crucial to avoid circular reasoning, where LLMs would be judging responses they themselves generated. This approach also helps mitigate potential bias.

To further investigate the influence of fixed randomness on LLM outputs, we systematically varied the temperature parameter during the LLM evaluation process. Temperature, a key hyperparameter in LLMs, significantly impacts the randomness (or "novelty") of their outputs (Peeperkorn et al., 2024; Davis et al., 2024). We evaluated LLM judgments at five temperature levels: 1.0, 0.75, 0.5, 0.25, and 0 (machine epsilon), while keeping the top-k sampling parameter constant at 50. This systematic variation in temperature allowed us to observe how changes in the level of randomness within the LLM's generation process affect the variability and reliability of the subsequent judgments, providing a window into the impact of fixed randomness.

The calculation of McDonald's omega employs the reliabilipy package (Fernández, 2022). It is important to note that in calculating the reliability, all judgments that did not provide a judgment of best response were treated as belonging to the same category ("Non Response") which was included as a category in addition to the five response categories "A" through "E". Consequently, the resulting calculated reliability serves as an upper bound on reliability because the variability in non-response judgments is ignored. This is further explored in the experiment results. Results are obtained using a pair of Quadro RTX 8000 GPUs with 48GB of memory each and the transformers library (Wolf et al., 2020).

6 Experimental Results

Our findings reveal several key insights. First, we observed that the ideal temperature for reliable LLM judgment is not universal; it varies depending on the specific LLM and the task, as revealed in Table 2. For example, some models (like Gemma-

³As of November 2024 at <https://huggingface.co/models>

Temperature	BBH	SQuAD	MT-Bench
1	0.702	0.632	0.462
0.75	0.713	0.639	0.618
0.5	0.703	0.644	0.476
0.25	0.677	0.598	0.64
0	1	1	1

(a) Omega reliability of Starling-LM-7B-beta judgments.

Temperature	BBH	SQuAD	MT-Bench
1	0.712	0.632	0.59
0.75	0.698	0.655	0.556
0.5	0.680	0.554	0.602
0.25	0.661	0.533	0.421
0	1	1	1

(b) Omega reliability of Meta-Llama-3-8B-Instruct judgments.

Temperature	BBH	SQuAD	MT-Bench
1	0.723	0.64	0.605
0.75	0.77	0.73	0.585
0.5	0.788	0.751	0.732
0.25	0.803	0.77	0.637
0	1	1	1

(c) Omega reliability of Gemma-1.1-7b-it judgments.

Table 2: Omega reliability by LLM judge across temperatures and benchmarks.

1.1-7b-it on SQuAD) exhibit decreased reliability at higher temperatures, while others (like Meta-Llama-3-8B-Instruct on BBH) show the opposite trend. Consequently, careful temperature tuning is crucial for each LLM and task to balance output novelty and judgment reliability. The ability to maintain reliability while increasing temperature, however, is a promising avenue for future research, as it provides opportunities to explore methods that generate novel, reliable, and high-quality outputs.

Reliability	Interpretation
$\alpha > 0.9$	Excellent
$0.9 > \alpha > 0.8$	Good
$0.8 > \alpha > 0.7$	Acceptable
$0.7 > \alpha > 0.6$	Questionable
$0.6 > \alpha > 0.5$	Poor
$\alpha < 0.5$	Unacceptable

Table 3: Rule of thumb for interpreting reliability measures.

Furthermore, the results in Table 2 highlight that

	Min	Q1	Q2	Q3	Max
BBH	0	0	0	1	2
SQuAD	3	5	5	5	5

Table 4: Five-number summary of accuracies of responses [A]-[E] across all sampled questions for each benchmark.

LLM judgment reliability is a major concern. It’s important to note that the calculated reliability values are likely optimistic estimates of the true reliability, as judgment outputs that did not select a response were treated as belonging to the same category (‘Non-response’). Despite this, the overall reliability scores are troubling. Adhering to the widely accepted reliability rule of thumb (Table 3), Starling-LM-7B-beta and Meta-Llama-3-8B-Instruct exhibited questionable reliability across benchmarks, with SQuAD and MT-Bench results consistently falling below the acceptable range (Tables 2a and 2b). This raises serious questions about the trustworthiness of their judgments.

Moreover, our results suggest a potential performance-reliability trade-off. The ranking of LLM judgment reliability (Gemma-1.1-7b-it > Starling-LM-7B-beta > Meta-Llama-3-8B-Instruct) inverts their Chatbot Arena performance ranking. This suggests that models optimized for performance may sacrifice reliability in their evaluations.

Finally, we find that benchmark quality impacts judgment reliability. SQuAD, despite its higher overall answer accuracy (Table 4), paradoxically exhibits lower judgment reliability than BBH (Table 2). We attribute this to the difficulty in distinguishing the best answer among multiple highly accurate responses in SQuAD. This highlights a critical weakness of current LLM judges: struggling with nuanced distinctions between correct answers.

7 Application

We explore the influence of LLM reliability on practical applications, focusing on how LLMs function as judges within the Head-to-Tail benchmark (Sun et al., 2023). This benchmark assesses LLM knowledge by posing questions about entities of varying popularity (“head,” “torso,” “tail”). We replicated the benchmark’s “head” question set for the Academic domain using Vicuna-7B, employing Starling-LM-7B-beta, Gemma-1.1-7b-it, and Meta-Llama-3-8B-Instruct as judges. Our goal is

Judge	A_{LM} Estimate
Sun et al., 2024	0.039 (N/A)
Starling-LM-7B-beta	0.0396 (0.00037)
Meta-Llama-3-8B-Instruct	0.0395 (0.00036)
Gemma-1.1-7b-it	0.0395 (0.00033)

Table 5: Average A_{LM} of Vicuna-7B responses to Head questions across replications of judgment models as compared to the results reported by Sun et al., 2024. Corresponding standard errors are provided beneath respective estimates for our judgment models.

Omega	Vicuna-7B
Starling-LM-7B-beta	0.9901
Meta-Llama-3-8B-Instruct	0.9896
Gemma-1.1-7b-it	0.9893

Table 6: Omega reliability of our judgment models Starling-LM-7B-beta, Gemma-1.1-7b-it, and Meta-Llama-3-8B-Instruct across replications for Head question responses by Vicuna-7B.

to demonstrate the strong connection between high reliability and consistent evaluation.

We replicated the A_{LM} metric (accuracy judged by an LLM) reported in Sun et al. (2023) using Vicuna-7B. Table 5 shows close agreement between our replicated A_{LM} values and the original results, confirming the reliability and reproducibility of LLM-based evaluations in this structured task. The consistent performance across different LLM judges further supports this.

Critically, our analysis reveals a strong link between high reliability and stable evaluation. Table 6 shows remarkably high reliability scores for our LLM judges across Head-to-Tail replications. Similarly, it is clear from the standard errors in the A_{LM} estimates in Table 5 that minimal variability exists across replications. This demonstrates that highly reliable LLM judges produce consistent and trustworthy evaluations, further supporting the validity of reliability.

It is unsurprising that the Head-to-Tail reliability is much higher than that of our experiments given that the LLM judgments in the Head-to-Tail paradigm only have two options ('correct' or 'incorrect') and are provided with a ground truth in the prompt. Thus, this employment of LLM-as-a-judge is much more simplistic than our

experiments and would be expected to show a higher reliability. While Head-to-Tail uses a simplified paradigm, it offers valuable insight into the reliability-consistency relationship. The observed high reliability, coupled with minimal variation in A_{LM} , strongly suggests that with careful design and clear criteria, LLMs can be dependable judges. Future research should extend these findings to more complex scenarios, developing methods to structure complex judgment tasks to leverage LLM reliability and enhance trustworthiness in diverse applications. Our findings reinforce reliability's importance in trustworthy LLM evaluation and highlight the potential for reliable LLM judgment in well-defined tasks.

8 Conclusions

Current LLM-as-a-judge methods, relying on single outputs, mask inherent judgment variability, creating a false sense of reliability (Ouyang et al., 2023; Wei et al., 2024). This "fixed randomness" (Section 1), poses risks, especially in high-stakes applications like medical AI evaluation. We address this by introducing a novel framework for evaluating LLM-as-a-judge reliability, focusing on internal consistency across replicated evaluations with varying random seeds — a crucial, previously overlooked aspect.

Our results validate our framework's importance. LLM judgment reliability is a significant concern, with several models exhibiting questionable reliability across benchmarks like SQuAD and MT-Bench. SQuAD's low reliability, despite high accuracy, reveals LLMs' difficulty with nuanced distinctions. Our research also suggests a performance-reliability trade-off, which our framework helps navigate. Finally, Head-to-Tail results further demonstrate our framework's value, as high reliability scores and minimal variability highlight the link between reliability and consistent evaluation.

This work provides a robust framework for assessing LLM-as-a-judge reliability, addressing a critical gap and developing a crucial tool for building trust in LLM-powered systems. As LLMs are increasingly used in complex applications, reliably assessing their judgment is paramount. Our framework empowers informed decisions about LLM deployment, promoting responsible usage. Prioritizing reliability is key to unlocking LLMs' potential while ensuring their ethical use.

9 Limitations

While this work represents a significant step forward in understanding and assessing the reliability of LLMs, it is important to acknowledge its limitations. Our work is currently focused on single and multi-turn benchmarks. These benchmarks were selected to provide a robust and generalizable understanding of the impact of reliability under varying and widely used contexts, including both structured and open-ended tasks. Future research can extend our approach to additional datasets and benchmarks. Further, developing domain-specific practical guidelines for interpreting and applying reliability metrics is an important area for future work given that different domains have varying tolerances for variability. For instance, fields like advertising may be more tolerant of variability, while fields like medicine require much higher levels of certainty. By providing domain-specific guidance, we can ensure that LLMs are used responsibly and effectively in various applications.

Despite these limitations, our work provides a strong foundation for future research on LLM reliability. By introducing a rigorous framework for evaluating LLM-as-a-judge responses, we offer valuable insights into the impact of reliability of LLM outputs and provide a roadmap for improving the quality and trustworthiness of LLMs.

10 Ethics Statement

This research does not raise any ethical concerns based on the theories or datasets employed. The benchmarks (BBH, SQuAD, MT-Bench, and Head-to-Tail) are publicly available and utilize publicly available data sources. As such, no privacy violations are anticipated. However, researchers should exercise caution when applying this work to proprietary datasets, as these may involve specific privacy regulations and ethical considerations.

References

- Aaron A Agbo. 2010. Cronbach’s alpha: Review of limitations and associated recommendations. *Journal of Psychology in Africa*, 20(2):233–239.
- AI@Meta. 2024. [Llama 3 model card](#).
- Sivakumar Alagumalai and David D Curtis. 2005. *Classical test theory*. Springer.
- Ebtessam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru,

Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.

Jacopo Amidei, Paul Piwek, and Alistair Willis. 2019. [Agreement is overrated: A plea for correlation to assess human evaluation reliability](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 344–354.

Berk Atil, Alexa Chittams, Liseng Fu, Ferhan Ture, Lixin Xu, and Breck Baldwin. 2024. [Llm stability: A detailed analysis with some surprises](#). *arXiv preprint arXiv:2408.04667*.

Shreya Bhandari, Yunting Liu, Yerin Kwak, and Zachary A Pardos. 2024. [Evaluating the psychometric properties of chatgpt-generated questions](#). *Computers and Education: Artificial Intelligence*, 7:100284.

Plaban Kumar Bhowmick, Anupam Basu, and Pabitra Mitra. 2008. [An agreement measure for determining inter-annotator reliability of human judgements on affective text](#). In *Coling 2008: Proceedings of the workshop on Human Judgements in Computational Linguistics*, pages 58–65.

Sairavi Kiran Biri, Subir Kumar, Muralidhar Panigrahi, Shaikat Mondal, Joshil Kumar Behera, and Himel Mondal. 2023. [Assessing the utilization of large language models in medical education: Insights from undergraduate medical students](#). *Cureus*, 15(10).

Chi-Min Chan, Yuxin Liu, Duyu Tang, Xiaohan Zhang, and Ming Zhang. 2023a. [Chateval: Towards better llm-based evaluators through multi-agent debate](#). *arXiv preprint arXiv:2308.07201*.

Colleen Chan, Kisung You, Sunny Chung, Mauro Giuffrè, Theo Saarinen, Niroop Rajashekhar, Yuan Pu, Yeo Eun Shin, Loren Laine, Ambrose Wong, et al. 2023b. [Assessing the usability of gutgpt: A simulation study of an ai clinical decision support system for gastrointestinal bleeding risk](#). *arXiv preprint arXiv:2312.10072*.

Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024. [Humans or llms as the judge? a study on judgement biases](#). *arXiv preprint arXiv:2402.10669*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#).

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, et al. 2024. [Chatbot arena: An open platform for evaluating llms by human preference](#). *arXiv preprint arXiv:2403.04132*.

- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. 2023. [A survey of chain of thought reasoning: Advances, frontiers and future](#). *arXiv preprint arXiv:2309.15402*.
- Jason Chuang, Margaret E Roberts, Brandon M Stewart, Rebecca Weiss, Dustin Tingley, Justin Grimmer, and Jeffrey Heer. 2015. [Topiccheck: Interactive alignment for assessing topic model stability](#). In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 175–184.
- Lee J Cronbach. 1951. [Coefficient alpha and the internal structure of tests](#). *psychometrika*, 16(3):297–334.
- Joshua Davis, Liesbet Van Bulck, Brigitte N Durieux, Charlotta Lindvall, et al. 2024. [The temperature feature of chatgpt: modifying creativity for clinical research](#). *JMIR Human Factors*, 11(1):e53559.
- Michael Desmond, Zahra Ashktorab, Qian Pan, Casey Dugan, and James M Johnson. 2024. [Evalullm: Llm assisted evaluation of generative outputs](#). In *Companion Proceedings of the 29th International Conference on Intelligent User Interfaces*, pages 30–32.
- Yijiang River Dong, Tiancheng Hu, and Nigel Collier. 2024. [Can llm be a personalized judge?](#) *arXiv preprint arXiv:2406.11657*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.
- Guohao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. 2024. [Towards revealing the mystery behind chain of thought: a theoretical perspective](#). *Advances in Neural Information Processing Systems*, 36.
- Rafael Valero Fernández. 2022. [reliability: measures of survey domain reliability in python with explanations and examples](#). *cronbach's alpha and omegas*. *Zenodo*.
- Xue-Yong Fu, Md Tahmid Rahman Laskar, Cheng Chen, and Shashi Bhushan TN. 2023a. [Are large language models reliable judges? a study on the factuality evaluation capabilities of llms](#). *arXiv preprint arXiv:2311.00681*.
- Xue-Yong Fu, Md Tahmid Rahman Laskar, Cheng Chen, and Shashi Bhushan Tn. 2023b. [Are large language models reliable judges? a study on the factuality evaluation capabilities of LLMs](#). In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 310–316, Singapore. Association for Computational Linguistics.
- Priyanshu Govil, Hemang Jain, Vamshi Krishna Bonagiri, Aman Chadha, Ponnurangam Kumaraguru, Manas Gaur, and Sanorita Dey. 2024. [Cobias: Contextual reliability in bias assessment](#). *arXiv preprint arXiv:2402.14889*.
- Roberto Gozalo-Brizuela and Eduardo C Garrido-Merchán. 2023. [A survey of generative ai applications](#). *arXiv preprint arXiv:2306.02781*.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. 2024. [A survey on llm-as-a-judge](#). *arXiv preprint arXiv:2411.15594*.
- Akash Gupta, Ivaxi Sheth, Vyas Raina, Mark Gales, and Mario Fritz. 2024. [Llm task interference: An initial study on the impact of task-switch in conversational history](#). *arXiv preprint arXiv:2402.18216*.
- Kevin A Hallgren. 2012. [Computing inter-rater reliability for observational data: an overview and tutorial](#). *Tutorials in quantitative methods for psychology*, 8(1):23.
- Johannes Hellrich and Udo Hahn. 2016. [Bad company—neighborhoods in neural embedding spaces considered harmful](#). In *Proceedings of coling 2016, the 26th international conference on computational linguistics: Technical papers*, pages 2785–2796.
- Robin K Henson. 2001. [Understanding internal consistency reliability estimates: A conceptual primer on coefficient alpha](#). *Measurement and evaluation in counseling and development*, 34(3):177–189.
- Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sébastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. 2023. [Phi-2: The surprising power of small language models](#). *Microsoft Research Blog*.
- Jaehun Jung, Faeze Brahman, and Yejin Choi. 2024. [Trust or escalate: Llm judges with provable guarantees for human agreement](#). *arXiv preprint arXiv:2407.18370*.
- Lisa Kollitsch, Klaus Eredics, Martin Marszałek, Michael Rauchenwald, Sabine D Brookman-May, Maximilian Burger, Katharina Körner-Riffard, and Matthias May. 2024. [How does artificial intelligence master urological board examinations? a comparative analysis of different large language models' accuracy and reliability in the 2022 in-service assessment of the european board of urology](#). *World Journal of Urology*, 42(1):20.
- Maia Kotelanski, Robert Gallo, Ashwin Nayak, and Thomas Savage. 2023. [Methods to estimate large language model confidence](#). *arXiv preprint arXiv:2312.03733*.
- Yusheng Liao, Yutong Meng, Hongcheng Liu, Yanfeng Wang, and Yu Wang. 2023. [An automatic evaluation](#)

- framework for multi-turn medical consultations capabilities of large language models. *arXiv preprint arXiv:2309.02077*.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*.
- Yan Liu, Amery D Wu, and Bruno D Zumbo. 2010. The impact of outliers on cronbach’s coefficient alpha estimate of reliability: Ordinal/rating scale item responses. *Educational and Psychological Measurement*, 70(1):5–21.
- Srija Macherla, Man Luo, Mihir Parmar, and Chitta Baral. 2023. Mddial: A multi-turn differential diagnosis dialogue dataset with reliability evaluation. *arXiv preprint arXiv:2308.08147*.
- Camila Paola Malkewitz, Philipp Schwall, Christian Meesters, and Jochen Hardt. 2023. Estimating reliability: A comparison of cronbach’s α , mcdonald’s ω_t and the greatest lower bound. *Social Sciences & Humanities Open*, 7(1):100368.
- Roderick P McDonald. 2013. *Test theory: A unified treatment*. psychology press.
- Brent Mittelstadt, Sandra Wachter, and Chris Russell. 2023. To protect science, we must use llms as zero-shot translators. *Nature Human Behaviour*, 7(11):1830–1832.
- Stefanie Nowak and Stefan Rüger. 2010. How reliable are annotations via crowdsourcing: a study about inter-annotator agreement for multi-label image annotation. In *Proceedings of the international conference on Multimedia information retrieval*, pages 557–566.
- Shuyin Ouyang, Jie M Zhang, Mark Harman, and Meng Wang. 2023. Llm is like a box of chocolates: the non-determinism of chatgpt in code generation. *arXiv preprint arXiv:2308.02828*.
- Bhrij Patel, Souradip Chakraborty, Wesley A Suttle, Mengdi Wang, Amrit Singh Bedi, and Dinesh Manocha. 2024. Aime: Ai system optimization via multiple lilm evaluators. *arXiv preprint arXiv:2410.03131*.
- Max Peepkorn, Tom Kouwenhoven, Dan Brown, and Anna Jordanous. 2024. Is temperature the creativity parameter of large language models? *arXiv preprint arXiv:2405.00492*.
- Harsh Raj, Vipul Gupta, Domenic Rosati, and Subhabrata Majumdar. 2023. Semantic consistency for assuring reliability of large language models. *arXiv preprint arXiv:2308.09138*.
- P Rajpurkar. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Jonas Rieger, Carsten Jentsch, and Jörg Rahnenführer. 2024. Ldaprototype: A model selection algorithm to improve reliability of latent dirichlet allocation. *PeerJ Computer Science*, 10:e2279.
- Kayla Schroeder and Zach Wood-Doughty. 2024. Reliability of topic modeling. *arXiv preprint arXiv:2410.23186*.
- Ruoxi Shang, Gary Hsieh, and Chirag Shah. 2024. Trusting your ai agent emotionally and cognitively: Development and validation of a semantic differential scale for ai trust. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pages 1343–1356.
- Xiaoming Shi, Jie Xu, Jinru Ding, Jiali Pang, Sichen Liu, Shuqing Luo, Xingwei Peng, Lu Lu, Haihong Yang, Mingtao Hu, et al. 2023. Llm-mini-cex: Automatic evaluation of large language model for diagnostic conversation. *arXiv preprint arXiv:2308.07635*.
- Claudio Spiess, David Gros, Kunal Suresh Pai, Michael Pradel, Md Rafiqul Islam Rabin, Amin Alipour, Susmit Jha, Prem Devanbu, and Toufique Ahmed. 2024. Calibration and correctness of language models for code. *arXiv preprint arXiv:2402.02047*.
- Kenneth Stensen and Stian Lydersen. 2022. Internal consistency: from alpha to omega. *Tidsskrift for den Norske Laegeforening: Tidsskrift for Praktisk Medicin, ny Raekke*, 142(12).
- Kai Sun, Yifan Ethan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. 2023. Head-to-tail: How knowledgeable are large language models (llm)? aka will llms replace knowledge graphs? *arXiv preprint arXiv:2308.10168*.
- Yuchong Sun, Che Liu, Kun Zhou, Jinwen Huang, Ruihua Song, Wayne Xin Zhao, Fuzheng Zhang, Di Zhang, and Kun Gai. 2024. Parrot: Enhancing multi-turn instruction following for large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9729–9750.
- Mirac Suzgun, Nathan Scales, Nathanael Schärfli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.
- Cheng Tan, Dongxin Lyu, Siyuan Li, Zhangyang Gao, Jingxuan Wei, Siqi Ma, Zicheng Liu, and Stan Z Li. 2024. Peer review as a multi-turn and long-context dialogue with role-based interactions. *arXiv preprint arXiv:2406.05688*.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024a. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, and et al. 2024b. [Gemma: Open models based on gemini research and technology](#). *Preprint*, arXiv:2403.08295.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). *arXiv preprint arXiv:2305.14975*.
- A Vaswani. 2017. [Attention is all you need](#). *Advances in Neural Information Processing Systems*.
- Sagar K Wadkar, Khajan Singh, Ritu Chakravarty, and Shivaji D Argade. 2016. [Assessing the reliability of attitude scale by cronbach's alpha](#). *Journal of Global Communication*, 9(2):113–117.
- Li Wang, Xi Chen, XiangWen Deng, Hao Wen, MingKe You, WeiZhi Liu, Qi Li, and Jian Li. 2024. [Prompt engineering in consistency and reliability with the evidence-based guideline for llms](#). *npj Digital Medicine*, 7(1):41.
- Hui Wei, Shenghua He, Tian Xia, Andy Wong, Jingyang Lin, and Mei Han. 2024. [Systematic evaluation of llm-as-a-judge in llm alignment tasks: Explainable metrics and diverse prompt templates](#). *arXiv preprint arXiv:2408.13006*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). *Advances in neural information processing systems*, 35:24824–24837.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pieric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.
- Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, et al. 2024. [Justice or prejudice? quantifying biases in llm-as-a-judge](#). *arXiv preprint arXiv:2410.02736*.
- Yiqun Zhang, Xiaocui Yang, Shi Feng, Daling Wang, Yifei Zhang, and Kaisong Song. 2024. [Can llms beat humans in debating? a dynamic multi-agent framework for competitive debate](#). *arXiv preprint arXiv:2408.04472*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Zangwei Zheng, Xiaozhe Ren, Fuzhao Xue, Yang Luo, Xin Jiang, and Yang You. 2024. [Response length perception and sequence scheduling: An llm-powered llm inference pipeline](#). *Advances in Neural Information Processing Systems*, 36.
- Ruiyang Zhou, Lu Chen, and Kai Yu. 2024. [Is llm a reliable reviewer? a comprehensive evaluation of llm on automatic paper reviewing tasks](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9340–9351.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, Karthik Ganeshan, Wei-Lin Chiang, Jian Zhang, and Jiantao Jiao. 2023. [Starling-7b: Improving llm helpfulness & harmlessness with rlaif](#).