

# DAFE: LLM-Based Evaluation Through Dynamic Arbitration for Free-Form Question-Answering

Anonymous ACL submission

## Abstract

Evaluating Large Language Models (LLMs) free-form generated responses remains a challenge due to their diverse and open-ended nature. Traditional supervised signal-based automatic metrics fail to capture semantic equivalence or handle the variability of open-ended responses, while human evaluation, though reliable, is resource-intensive. Leveraging LLMs as evaluators offers a promising alternative due to their strong language understanding and instruction-following capabilities. Taking advantage of these capabilities, we propose the Dynamic Arbitration Framework for Evaluation (DAFE), which employs two primary LLM-as-judges and engages a third arbitrator only in cases of disagreements. This selective arbitration prioritizes evaluation reliability while reducing unnecessary computational demands compared to conventional majority voting. DAFE utilizes task-specific reference answers with dynamic arbitration to enhance judgment accuracy, resulting in significant improvements in evaluation metrics such as Macro F1 and Cohen’s Kappa. Through experiments, including a comprehensive human evaluation, we demonstrate DAFE’s ability to provide consistent, scalable, and resource-efficient assessments, establishing it as a robust framework for evaluating free-form model outputs.

## 1 Introduction

The rapid advancements in Large Language Models (LLMs) have propelled the field of natural language processing forward, yet their evaluation remains a challenge (Laskar et al., 2024). In particular, free-form model responses are difficult to evaluate because their correctness depends on understanding the broader context and underlying meaning (Si et al., 2021). Many benchmarks, such as MMLU (Hendrycks et al., 2021), often simplify evaluation by focusing on structured formats (e.g., multiple-choice questions) (Chen et al., 2024). Although effective for certain tasks, such methods

rely on log probabilities assigned to predefined options, where the model selects the most likely answer, limiting the range of capabilities that can be assessed (Thakur et al., 2024). This structured approach fails to accommodate the complexity of free-form responses, where multiple valid answers exist (Chang et al., 2024). The rigid, predefined options in such evaluations not only limit the scope of assessment but also overlook the diversity of potential correct responses in free-form tasks (Li et al., 2023; Zhang et al., 2024).

Automatic metrics including lexical matching, n-gram, and neural-based have been widely adopted as scalable solutions for the evaluation of free-form model outputs. Lexical matching methods such as Exact Match (EM) evaluate model predictions by assessing strict lexical alignment between generated outputs and reference answers. However, EM fails to account for semantically equivalent variations in phrasing. For instance, despite their equivalence, EM treats “nuclear weapon” and “atomic bomb” as incorrect. Similarly, n-gram-based metrics (Papineni et al., 2002; Lin, 2004) primarily assess surface-level similarity and often fail to capture semantic equivalence, particularly when lexical or structural diversity conveys the same underlying meaning (Zhu et al., 2023; Chen et al., 2021; Zhang et al., 2020). Neural-based metrics like BERTScore (Zhang et al., 2020) address such limitations by leveraging contextual embeddings to evaluate semantic similarity. However, BERTScore depends on reference quality (Liu et al., 2024) and struggles with domain adaptation and length variations (Zhu et al., 2023). Furthermore, continuous score provider metrics are difficult to interpret (Xu et al., 2023). The limitations in automatic metrics become particularly evident when evaluating instruction-tuned chat models (Doostmohammadi et al., 2024), which tend to produce verbose and diverse responses (Saito et al., 2023; Wang et al., 2024b).

Contrary to automatic metrics, human evaluation provides a more transparent assessment (Chiang and Lee, 2023). However, despite being the “gold standard”, human evaluation is not without its limitations. LLMs’ growing complexity and scale have made recruiting and coordinating multiple human raters increasingly resource-intensive and time-consuming (Mañas et al., 2024). Furthermore, the reliability of human evaluation is additionally challenged by variations in rater expertise and inherent subjectivity that affect reproducibility (Clark et al., 2021; Chiang and Lee, 2023).

Recently, a paradigm shift has emerged where LLMs are utilized to judge the candidate model generations for given tasks (Zheng et al., 2024). This model-based method leverages the instruction-following capabilities of LLMs through evaluation prompts or, in some cases, fine-tuned versions of LLMs that are specifically optimized for evaluation. In this new line of work, research primarily focuses on pairwise comparison (Zheng et al., 2024; Wang et al., 2023; Vu et al., 2024), such as instructing an LLM to judge “which assistant response is better”, and single-answer scoring (Verga et al., 2024) like evaluating summarization task based on predefined criteria (e.g., likability, relevance, etc.) (Chiang and Lee, 2023; Hu et al., 2024; Liu et al., 2023; Chan et al., 2024; Chu et al., 2024).

Inspired by a recent study on self-correction where external feedback helps models identify and correct their mistakes (Gou et al., 2024a), we propose to guide LLM-as-a-judge with human-annotated task-specific reference answers in order to explore the potential of LLMs as an alternative to lexical matching (e.g., EM), neural-based (e.g., BERTScore), and human evaluation for automatic evaluation of free-form model responses. Unlike traditional metrics, an LLM judge can leverage its language understanding and instruction-following capabilities to recognize the correctness of open-ended generations.

We propose the Dynamic Arbitration Framework for Evaluation (DAFE), which employs LLM judges to evaluate free-form model responses. Using a single LLM as a judge, while simple, often leads to inconsistent evaluations, undermining trust in the results. On the other hand, the common practice of using large, universally capable models such as GPT-4 as evaluators makes the evaluation process both slow and costly (Jung et al., 2024; Adlakha et al., 2024; Verga et al., 2024), further limit-

ing its broader applicability. Relying on multiple judges for every evaluation, though more reliable, exacerbates these computational challenges, making such approaches impractical at scale. DAFE offers a middle ground between these approaches by utilizing two complementary primary judges to perform the initial assessment. Only when these judges disagree, is a third independent arbitrator engaged to resolve the conflict. This selective arbitration ensures evaluation reliability and fairness while reducing computational overhead. Our experiments reveal that DAFE achieves significant improvements in metrics such as Macro F1 and Cohen’s kappa. Our key contributions include: a detailed analysis of limitations in conventional metrics for free-form QA, an evaluation of LLM judges with insights into their strengths and errors, a comprehensive human evaluation for benchmarking, and the introduction of DAFE—a scalable framework that improves reliability while minimizing the need for additional evaluators through selective arbitration.

## 2 Methodology

This section briefly describes the key components of our proposed framework.

### 2.1 Candidate LLMs

A candidate LLM  $\mathcal{C}_{\text{llm}}$  generates output  $\bar{y}$  for the given input  $x$ . We first utilized candidate LLMs to obtain outputs for the given free-form question-answering tasks.

### 2.2 LLMs-as-a-Judge

A judge  $\mathcal{J}_{\text{llm}}$  LLM delivers evaluation or verdict  $V$  on candidate LLMs  $\mathcal{C}_{\text{llm}}$  outputs  $\bar{y}$ . The  $\mathcal{J}_{\text{llm}}$  evaluates output when prompted with  $x$  (i.e.,  $x \rightarrow \mathcal{A}_{\text{llm}}$ ) and  $\bar{y}$ . We utilized the reference answer  $r$  and prompted  $P$  the  $\mathcal{J}_{\text{llm}}$  as:

$$P = \{x, \bar{y}, r\}$$

Utilizing  $P$ ,  $\mathcal{J}_{\text{llm}}$  performs the evaluation and delivers a decision as  $V = J(P)$ . The structure of this  $V$  depends on the instructions provided in  $P$ . For instance, if a binary  $V$  is required,  $J$  assesses whether  $\bar{y}$  is aligned with  $r$  given the context  $x$  and returns True if  $\bar{y}$  is deemed correct, or False if it is not. The evaluation  $P$  may vary from zero-shot, where  $\mathcal{J}_{\text{llm}}$  receives no prior examples, to few-shot, which includes several related examples, or a chain of thought, encouraging  $\mathcal{J}_{\text{llm}}$  to reason stepwise through the problem.

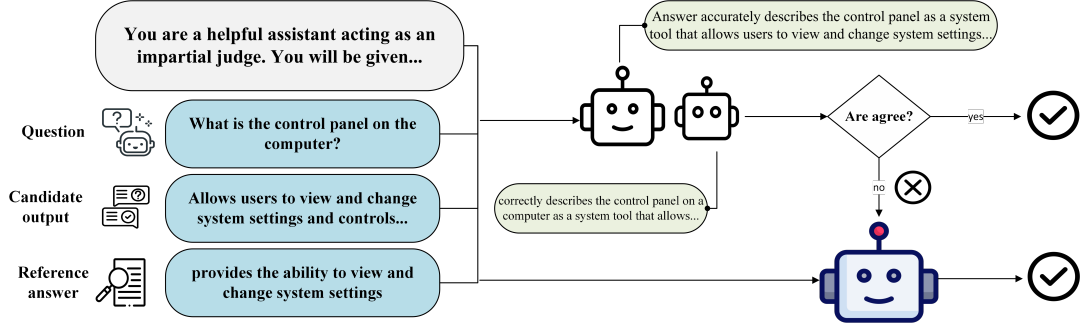


Figure 1: **Our proposed Dynamic Arbitration Framework for Evaluation (DAFE)**. Two primary judges,  $J_1$  and  $J_2$ , first provide verdicts  $V_{i_1}$  and  $V_{i_2}$  for an instance  $i$ . If agree, that consensus  $V_i$  is the final decision  $D_i$ . If disagree, a tiebreaker model  $J_t$  independently produces a verdict  $V_t$ . The final decision  $D_i$  is then determined via majority voting among  $\{V_{i_1}, V_{i_2}, V_t\}$ .

### 2.3 Dynamic Arbitration Framework for Evaluation (DAFE)

In traditional human evaluation settings, when two annotators disagree on a judgment, a third expert is often called upon to resolve the dispute. Drawing inspiration from this efficient human arbitration practice, we propose the Dynamic Arbitration Framework for Evaluation (DAFE). Rather than immediately employing a large powerful or a closed-source LLMs-as-a-judge, DAFE adopts a cost-efficient approach by beginning with two complementary open-source models as primary judges based on their past performance (Kenton et al., 2024). When these judges reach a consensus, no further evaluation is needed. Only in cases of disagreement is the more powerful LLM engaged as an arbitrator, whose decision then creates a majority verdict. This dynamic approach maintains evaluation quality while minimizing reliance on expensive models. The method also accounts for varying skill levels across different LLMs and tasks (Liang et al., 2024; Sun et al., 2024).

Formally, let  $V_{i_1}$  and  $V_{i_2}$  denote the verdicts from the two primary judges for the  $i$ -th evaluation instance. We define the agreement status  $A_i$  as:

$$A_i = \begin{cases} 1 & \text{if } V_{i_1} = V_{i_2}, \\ 0 & \text{otherwise.} \end{cases}$$

If  $A_i = 1$ , the final decision  $D_i$  is simply  $V_i$ , the agreed-upon verdict of the primary judges. If  $A_i = 0$ , a tiebreaker model provides an additional verdict  $V_t$ . The final decision  $D_i$  is then obtained via majority voting among  $\{V_{i_1}, V_{i_2}, V_t\}$ .

Formally:

$$D_i = \begin{cases} V_i & \text{if } A_i = 1, \\ \text{majority}(\{V_{i_1}, V_{i_2}, V_t\}) & \text{if } A_i = 0. \end{cases}$$

The majority operation selects the verdict that appears at least twice among  $\{V_{i_1}, V_{i_2}, V_t\}$ . Since there are three votes, at least two must coincide for a majority.

## 3 Experiments

We utilize the following settings to examine the performance and reliability of individual LLM judges and DAFE.

### 3.1 Models

We select open and closed-source instruct models to serve as candidates and judges in our experiment. These include DeepSeek-V3 (DeepSeek-AI et al., 2025), Llama-3.1 70B<sup>1</sup> (Meta AI, 2024), GPT-3.5-turbo (Brown et al., 2020), Mistral 7B<sup>2</sup> (Jiang et al., 2023), and Mixtral 8x7B<sup>3</sup> (Jiang et al., 2024). We also utilize GPT-4o (OpenAI et al., 2023) and DeepSeek-R1 (DeepSeek-AI et al., 2025) in our ablation experiments. To ensure the reproducibility of our experiments, we set the temperature to 0 for all models under study, as the performance of LLM-based evaluators has been shown to drop when temperature increases (Hada et al., 2024). For our proposed DAFE method, we utilized Mistral 7B and Llama 3.1 70B as primary judges with

<sup>1</sup><https://huggingface.co/meta-llama/Meta-Llama-3.1-70B-Instruct>

<sup>2</sup><https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

<sup>3</sup><https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

GPT-3.5-turbo as the tiebreaker. In addition, we experiment with other models as tiebreakers in our ablation experiments. In the rest of the paper, we refer both candidate and judge LLMs as: DeepSeek, Llama, GPT, Mistral, and Mixtral.

### 3.2 Datasets

We focus on free-form question-answering (QA) since it has widespread practical applications and the critical importance of truthfulness in this domain (Gou et al., 2024a; Evans et al., 2021). In our experiment, we utilize five free-form QA datasets: AmbigQA (Min et al., 2020), FreshQA (Vu et al., 2023), HotpotQA (Yang et al., 2018), Natural Questions (Kwiatkowski et al., 2019), and TriviaQA (Joshi et al., 2017). See Appendix A for details.

### 3.3 Prompts

We designed generalized (i.e., with minimum instructions) zero-shot prompts with role-playing (Kong et al., 2024) for both candidates and judges. Initially, we prompt candidate LLMs to elicit outputs for the given random samples associated with each dataset.

To evaluate the outputs of candidate LLMs, we prompt judge LLMs for binary verdicts (i.e., True or False) using  $P = \{x, \bar{y}, r\}$  and instructed to provide a brief explanation for their verdicts (see Appendix D for examples). Binary verdicts explicitly differentiate between correct and incorrect answers, minimize subjective interpretations, and simplify the evaluation process, thus facilitating automatic evaluation. In addition to three key prompt components (i.e.,  $x, \bar{y}, r$ ), we define the role of the judge LLMs as “*You are a helpful assistant acting as an impartial judge.*” to mitigate biases in judgments (Zheng et al., 2024). We chose not to use few-shot or chain-of-thought prompting strategies to keep the solution robust to a variety of tasks. Previous studies have also shown that in-context examples do not significantly improve the performance of model-based evaluators (Hada et al., 2024; Min et al., 2022).

### 3.4 Baselines

We establish the following baselines.

**Exact Match (EM):** For our selected datasets and also free-form QA tasks, EM serves as a standard lexical matching metric to evaluate candidate LLM performance (Izacard and Grave, 2021; Lewis

et al., 2020; Gou et al., 2024b). Due to the verbose nature of LLM-generated responses, we adapt EM to classify an answer as correct if any golden answer  $r_i \in R$  appears within the generated response  $\bar{y}$  (i.e.,  $r_i \subseteq \bar{y}$ ), rather than requiring complete strict string equality (i.e.,  $\bar{y} = r_i$ ).

**BERTScore:** We use BERTScore (Zhang et al., 2020) which measures similarity by comparing contextualized word embeddings derived from a pre-trained BERT model. This enables the evaluation to focus on semantic correctness rather than exact lexical matches. As BERTScore is based on continuous values between -1 and 1, we set a threshold of  $\tau = 0.5$  to convert continuous similarity scores into binary 0 and 1. The purpose of this conversion is to allow direct comparison with other evaluation methods. For our implementation, we use the microsoft/deberta-xlarge-mnli<sup>4</sup> model (He et al., 2021).

**G-Eval:** In addition to automatic metrics, we also utilize G-Eval (Liu et al., 2023), a reference-free framework that uses GPT-4 to assess the quality of the generated text. In this setting, we modify the evaluation prompt by excluding the reference answer  $r$  and directly prompted the evaluator model as  $P = \{x, \bar{y}\}$  along with instructions.

**Human Evaluation:** It remains the gold standard for assessing the outputs of candidate LLMs. We recruit three graduate students from our academic network, all specialized in natural language processing, to serve as annotators. We provide the input given to the candidate LLMs, reference answers, and candidate LLMs responses. This format, while similar, is distinct from the judge LLMs prompts which additionally require formatted decisions. We anonymize the origin of model responses to reduce potential bias linked to model familiarity or reputation. The annotators were asked to score the candidate LLMs outputs on a binary scale: ‘1’ for ‘True’ and ‘0’ for ‘False’ based on alignment with the reference answer and contextual relevance. For inter-rater reliability, we compute Fleiss’ Kappa ( $\kappa$ ) (Fleiss and Cohen, 1973) and percent agreement. See Appendix B for details.

## 4 Results

Figure 2 illustrates the raw performance of Llama obtained through various evaluators. Unlike lexical

<sup>4</sup><https://huggingface.co/microsoft/deberta-xlarge-mnli>



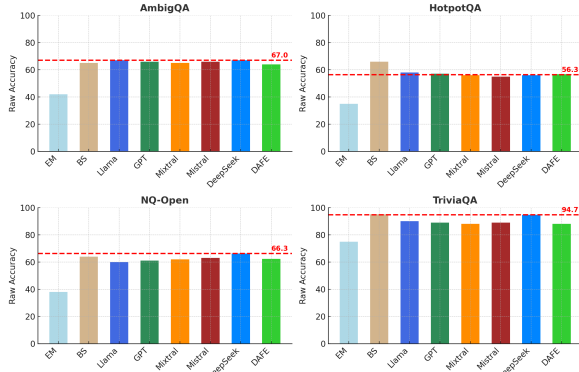


Figure 2: Raw accuracy of candidate Llama across free-form QA tasks using automatic metrics and model-based evaluation. The Human Majority (HM) serves as the ground truth for all evaluators.

matching and neural-based metrics, each LLM-as-a-judge shows overall performance close to the human majority. The proposed DAFE method consistently achieves comparable or slightly better alignment with the human majority. Conventional metrics such as EM severely underestimate the candidate LLMs’ performance. Contrarily, BERTScore tends to overestimate the performance except in some cases such as when evaluating Llama on AmbigQA and NQ-Open (see Table 6 in Appendix C for additional results).

#### 4.1 Alignment with human evaluation

We calculate Cohen’s kappa (McHugh, 2012) to find the agreement between each evaluator and the human majority to obtain instance-level comparison. Overall, DAFE is almost perfectly aligned with human judgment than other evaluators (see Table 1). Similarly, individual LLM judges show more substantial to nearly perfect agreement with human judgments than EM and BERTScore.

Due to the high-class imbalance in TriviaQA, kappa scores can be misleadingly low despite high raw agreement - a known limitation called the “kappa paradox” (Cicchetti and Feinstein, 1990). Therefore, we treat the evaluation as a binary classification task where we consider each evaluator’s predictions against the human majority and report Macro-F1 scores which give equal weight to both classes regardless of their frequency in the selected random samples.

As evidenced by consistently high Macro F1 scores in Table 2, DAFE maintains a strong alignment with human judgment. This represents a substantial improvement over individual model perfor-

Table 1: Cohen’s Kappa scores displaying the agreement levels of various evaluators with human judgments across candidate models and tasks. Higher scores indicate better agreement with human judgments.

LLMs	Tasks	Evaluators							
		EM	BS	DeepSeek	Llama	GPT	Mixtral	Mistral	DAFE
Llama	AmbigQA	0.518	0.283	0.897	0.888	0.844	0.824	0.858	0.911
	HotpotQA	0.577	0.498	0.885	0.877	0.899	0.820	0.832	0.953
	NQ-Open	0.381	0.437	0.797	0.833	0.793	0.816	0.738	0.927
	TriviaQA	0.281	0.564	0.460	0.547	0.439	0.396	0.299	0.684
GPT	AmbigQA	0.561	0.252	0.951	0.944	0.897	0.861	0.853	0.967
	HotpotQA	0.604	0.300	0.807	0.953	0.973	0.873	0.933	0.987
	NQ-Open	0.453	0.218	0.809	0.884	0.824	0.824	0.829	0.956
	TriviaQA	0.335	0.364	0.594	0.650	0.401	0.580	0.467	0.775
Mixtral	AmbigQA	0.546	0.337	0.896	0.896	0.781	0.909	0.887	0.951
	HotpotQA	0.546	0.349	0.920	0.940	0.933	0.859	0.940	0.973
	NQ-Open	0.371	0.301	0.825	0.879	0.728	0.899	0.815	0.913
	TriviaQA	0.317	0.390	0.661	0.625	0.605	0.678	0.436	0.764
Mistral	AmbigQA	0.599	0.254	0.893	0.893	0.893	0.893	0.860	0.953
	HotpotQA	0.605	0.383	0.903	0.937	0.902	0.895	0.937	0.958
	NQ-Open	0.484	0.291	0.797	0.851	0.838	0.878	0.840	0.953
	TriviaQA	0.467	0.239	0.754	0.758	0.725	0.645	0.470	0.854

mance, where individual judges generally revealed varying levels of agreement with human evaluation. LLM-as-a-judge approach generally works better with larger more powerful models. This is particularly noticeable in DeepSeek and GPT which achieve higher Macro-F1 scores (0.97-0.98) across AmbigQA, HotpotQA, and NQ-Open compared to smaller models. This reveals an important scaling law in evaluation capability (Kaplan et al., 2020; Zheng et al., 2024; OpenAI et al., 2024). However, we also found that the most advanced models are not always guaranteed to be the best evaluators. We observed slightly comparable performance through small open-source Mistral-7B. For instance, when evaluating candidate Mixtral-8x7B on AmbigQA, Mistral-7B as-a-judge outperformed (0.944) judge GPT-3.5-turbo (0.891). Regardless, we observe relatively lower Macro-F1 scores for all LLM judges in TriviaQA.

Interestingly, despite EM’s deviation from the human majority (see Figure 2 and Table 6), lexical matching EM typically accomplishes better alignment with human evaluation on instance-level in Table 2 than neural-based BERTScore. EM’s strict and conservative nature leads to lower overall performance, but its high-precision characteristics ensure that when it identifies a match, it strongly aligns with human judgment. In contrast, BERTScore takes a more lenient approach to semantic matching. Although this leniency produces higher raw scores, it introduces more false positives, consequently reducing instance-level agreement with human judgments. This pattern emerges clearly in many models and tasks such as when evaluating Llama-3.1-70B on AmbigQA, EM shows a raw score of 42.3% but achieves a Macro-F1 of

Table 2: Macro-F1 scores of various evaluators applied to different candidate LLMs and associated tasks. Higher scores indicate better performance. DAFE consistently achieves the highest Macro-F1 across all evaluated settings.

LLMs	Tasks	Evaluators							
		EM	BS	DeepSeek	Llama	GPT	Mixtral	Mistral	DAFE
Llama	AmbigQA	0.744	0.641	0.948	0.944	0.922	0.912	0.929	0.955
	HotpotQA	0.778	0.745	0.942	0.939	0.949	0.910	0.916	0.976
	NQ-Open	0.653	0.718	0.898	0.916	0.896	0.907	0.869	0.964
	TriviaQA	0.612	0.782	0.726	0.772	0.717	0.695	0.640	0.842
GPT	AmbigQA	0.792	0.622	0.976	0.972	0.949	0.930	0.927	0.984
	HotpotQA	0.794	0.623	0.903	0.977	0.987	0.936	0.966	0.993
	NQ-Open	0.703	0.606	0.904	0.942	0.911	0.911	0.914	0.978
	TriviaQA	0.646	0.681	0.796	0.824	0.700	0.789	0.730	0.887
Mixtral	AmbigQA	0.760	0.666	0.948	0.948	0.891	0.955	0.944	0.975
	HotpotQA	0.761	0.657	0.960	0.970	0.966	0.930	0.970	0.987
	NQ-Open	0.650	0.649	0.912	0.939	0.863	0.950	0.908	0.956
	TriviaQA	0.625	0.695	0.829	0.812	0.803	0.838	0.716	0.882
Mistral	AmbigQA	0.792	0.622	0.947	0.947	0.947	0.947	0.930	0.977
	HotpotQA	0.796	0.673	0.951	0.969	0.951	0.947	0.969	0.979
	NQ-Open	0.726	0.639	0.898	0.925	0.919	0.939	0.920	0.976
	TriviaQA	0.718	0.608	0.925	0.879	0.863	0.822	0.735	0.927

0.744, while BERTScore indicates a higher raw score of 63.0% but a lower Macro-F1 of 0.641.

## 4.2 Analysis

In our experiments, candidate LLMs generated 7,500 outputs for the given tasks, with each evaluator producing 7,500 corresponding evaluations. We randomly sampled 100 error cases (50 false positives and 50 false negatives) from each evaluator to understand their behavior. Given EM had only 11 false positives, we included all of them in our analysis. Due to space constraints, we moved the detailed analysis of EM and BERTScore to Appendix C and focused exclusively on the LLM-as-a-judge method here.

LLM-based evaluators demonstrate strong abilities in recognizing semantic variations while maintaining the core meaning, especially when assessing responses that use different terminology or structural approaches to convey the same information. For instance, in the evaluation examples, evaluators correctly identified that “Salma Hayek” and “Salma Hayek Pinault” refer to the same individual, acknowledging the semantic equivalence despite differences in phrasing. Similarly, when assessing responses that use different terms for the same entity, such as recognizing “Nick Fury, Agent of S.H.I.E.L.D.” as part of the broader “Marvel” universe, the evaluators effectively maintain the core meaning and contextual relevance. Their explanations show systematic assessment patterns that combine multiple evaluation criteria including factual accuracy, logical coherence, and contextual relevance.

LLMs are prone to hallucination in justifica-

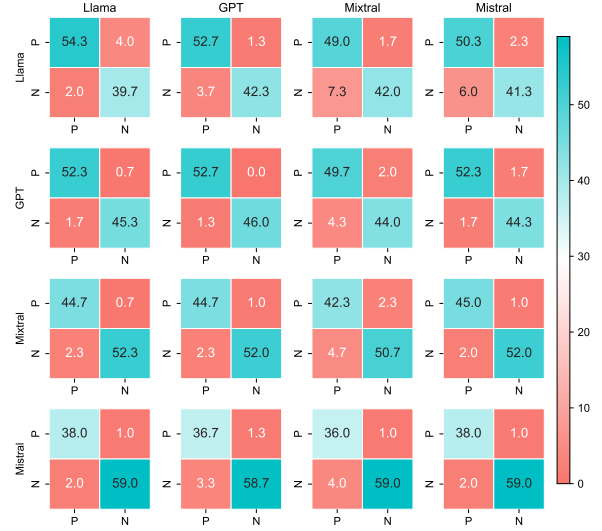


Figure 3: Heatmap illustrating the performance of four individual LLM judges on HotpotQA. Each cell value represents percentages (%). Rows represent predicted outcomes (P: Positive, N: Negative), while columns represent actual outcomes. See Appendix C for full results.

tion (Zhang et al., 2023), where they fabricate reasoning to support their evaluations, produce detailed but incorrect explanations, or reference non-existent criteria or standards. In LLM judges, false positives and negatives (e.g., see Figure 3) often result from overlooking critical distinctions between candidate LLM outputs and failing to account for the specificity required by the reference answer. This pattern is particularly noticeable in Mistral 7B, where the model disregards the ground truth and provides evaluations influenced by unknown factors. For example, when evaluating candidate GPT-3.5’s response “The foreign minister of Germany who signed the Treaty of Versailles was Hermann Müller.” which is correct according to the reference answer “Hermann Müller” and human evaluation, Mistral 7B as-a-judge incorrectly marked this response as false and fabricated reasoning “Hermann Müller was the Chancellor of Germany, not the Foreign Minister. The Foreign Minister of Germany who signed the Treaty of Versailles was Gustav Stresemann.” in support of its decision. The same problem can also be attributed to inconsistent evaluations. Because when Mistral 7B acted as a candidate for the same question, its response to the question is completely different: “The Treaty of Versailles was signed by Matthias Erzberger, a German politician who served as the President of the German National Assembly at the

time”. There are also alternative interpretations of this issue, such as ambiguity in the question, but we leave a deeper exploration of these aspects to future work.

We observe a different pattern in some judges, specifically, GPT-3.5 and Mixtral 8x7B which focuses more on specificity. This approach shifts the evaluation towards false negatives by missing semantically similar but structurally different answers. We found many cases when such evaluators failed to account for valid variations in phrasing or granularity, focusing instead on rigid adherence to the reference answer. Compounding these issues are reasoning errors within the evaluators’ own explanations, which often contain fabrications, circular logic, or overconfident assertions. By insisting on correctness derived strictly from the reference, evaluators disregard valid alternative perspectives and can even mischaracterize or invert the facts in their attempts to justify their decisions. This dynamic leaves little room for nuance or ambiguity, and it pushes the evaluation process away from fair, context-sensitive assessment toward rigid, and sometimes inaccurate, verdicts.

Verbosity (Ye et al., 2024) emerges as a subtle source of bias, where more elaborate answers are sometimes overrated simply due to their detail and fluency, while concise yet correct responses are undervalued. This misplaced emphasis leads to irrelevant judgment criteria, such as praising the presence of irrelevant information or penalizing perfectly valid but succinct answers. We also found that LLM-based judges encounter challenges in multiple reference answers and more open-ended questions. This confusion is especially pronounced in the TriviaQA where the diversity and flexibility of valid responses present challenges for the judges’ ability to consistently recognize and evaluate a range of correct answers.

We found several temporal limitations in LLM-based evaluators. Although most of our datasets are older and the evaluator models are relatively up-to-date, we still observed instances where references to recent events, newly emerging terminology, or evolving contexts were misinterpreted. The FreshQA dataset (Vu et al., 2023), being recent, serves as a valuable testbed for assessing these temporal deficiencies. As shown in Table 3, LLM-based evaluators indicate deviation from human judgment on FreshQA compared to tasks that rely on older information, such as HotpotQA. Specif-

Table 3: Performance (in Macro F1) of LLM judges on FreshQA.

LLMs	Evaluators					
	DeepSeek	Llama	GPT	Mixtral	Mistral	DAFE
DeepSeek	0.714	0.692	0.715	0.614	0.724	0.830
Llama	0.801	0.835	0.737	0.817	0.730	0.917
GPT	0.659	0.695	0.824	0.780	0.746	0.891
Mixtral	0.732	0.708	0.779	0.738	0.703	0.936
Mistral	0.687	0.665	0.802	0.818	0.723	0.880

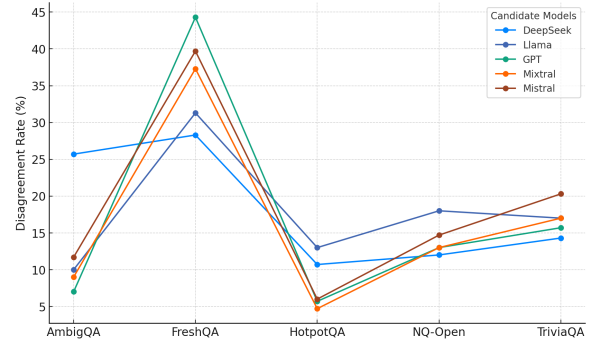


Figure 4: Disagreement rates between primary judges (Llama+Mistral) across candidate LLMs and tasks.

ically, in dynamic or time-sensitive contexts, we found that LLM judges tend to hallucinate by consistently classifying candidate model responses as True, even when incorrect. For example, when presented with the question: “On what date did the Patriots last play the Miami Dolphins?” the LLM-generated response states: “The last time the New England Patriots played the Miami Dolphins was on January 1, 2023, during the NFL regular season.” Despite the correct reference answer being “November 24, 2024” the LLM evaluator not only failed to recognize the inaccuracy but also hallucinated an erroneous justification, stating: “The proposed answer correctly states the date the New England Patriots last played the Miami Dolphins as January 1, 2023, which matches the information provided.”

### 4.3 Disagreements between primary judges

Figure 4 shows that disagreements between our primary judges, Llama-3.1 70B and Mistral 7B, mainly occur in the TriviaQA and FreshQA, with disagreement rates reaching 20.3% and 44.3%, respectively. Interestingly, higher disagreement rates between primary judges create a greater opportunity for DAFE to refine evaluations. As depicted in Figure 4, FreshQA (31.3% for Llama-70B, 39.7% for Mistral-7B) demonstrates the highest disagreement, allowing DAFE to improve Macro F1 scores



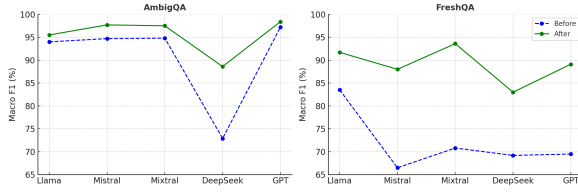


Figure 5: Comparison of Macro F1 scores before and after arbitration (see Appendix C for more results).

(see Table 3).

#### 4.4 Impact of arbitration

Our proposed arbitration approach significantly enhanced evaluation performance by resolving disputes through an independent judge, GPT-3.5-turbo (see Figure 5). Notably, in the AmbigQA, Macro F1 scores advanced from 72.9% to 86.6%, and Cohen’s Kappa increased from 0.467 to 0.773 (see Figure 7). These improvements highlight the pivotal role of the arbitrator in ensuring reliable and consistent evaluation outcomes.

### 5 Related work

Evaluation of natural language generation has traditionally relied on supervised signal-based metrics such as EM which evaluates the exact lexical match between generated outputs and reference answers. Despite its simplicity and efficiency, EM overlooks semantically equivalent variations, often penalizing accurate responses that use different phrasing (Wang et al., 2024a; Kamaloo et al., 2023). Other commonly used metrics including BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) primarily focus on n-gram overlap with human written reference texts. Despite their widespread use, these metrics have significant limitations in capturing semantic subtleties and contextual relevance (Zhang et al., 2020). To address the limitations of conventional metrics, various model-based methods such as BERTScore (Zhang et al., 2020) offer semantically informed evaluation. However, even BERTScore and similar embedding-based methods struggle to effectively evaluate open-ended generation (Zheng et al., 2024; Sun et al., 2022).

Recent advances in LLMs have unlocked new opportunities for automatic and context-aware evaluation (Li et al., 2024b; Chiang and Lee, 2023; Zheng et al., 2024). A key strength of LLM-based evaluators lies in their ability to operate in reference-free settings, where evaluation does not rely on pre-defined answers but instead leverages subjective

criteria such as helpfulness, relevance, and coherence. This capability makes LLM evaluators particularly well-suited for assessing tasks where multiple valid responses exist or where human-like judgment is required (Li et al., 2024a). For instance, LLMs are frequently used in subjective evaluations such as pairwise comparison (“Which response is better?”) or single-response scoring (“How good is this response based on criteria X?”) (Verga et al., 2024; Chan et al., 2024). LLM-based evaluators are specifically effective for tasks like summarization, where subjective criteria are central to evaluation (Liu et al., 2023). However, they are less effective for fact-based tasks such as free-form question-answering, where responses are either correct or incorrect and require explicit verification against reference answers.

Furthermore, LLM-based evaluators face several challenges, particularly in ensuring consistency and fairness (Ye et al., 2024; Khan et al., 2024). In reference-free settings, the absence of a definitive ground truth increases the risk of bias in evaluations (Ye et al., 2024; Kim et al., 2024; Huang et al., 2024a). Common biases include positional bias, where LLMs may favor responses based on their order (Zheng et al., 2024; Khan et al., 2024), verbosity bias, which favors longer or more detailed responses (Huang et al., 2024b), and self-enhancement bias, where models may disproportionately prefer their own outputs (Zheng et al., 2024). These biases can distort evaluations and undermine the reliability of the results.

### 6 Conclusion

We present DAFE, a framework designed to evaluate free-form question-answering by leveraging LLMs. Our findings demonstrate that individual LLM judges are reliable alternatives to traditional lexical and neural-based metrics, offering closer alignment with human evaluations. However, relying solely on individual judges poses challenges including inherent biases and prompt sensitivity, which can affect evaluation performance. DAFE addresses these challenges through a dynamic arbitration mechanism. This design achieves near-perfect agreement with human evaluations, establishing DAFE as a trustworthy and reliable framework for evaluating open-ended language generation tasks. In the future, we aim to explore DAFE by excluding reference answers and integrating LLM agents with tools-interacting capabilities for evaluation.



## 7 Limitations

We acknowledge certain limitations in our study. The accuracy of evaluations depends on the quality and clarity of reference answers, which serve as the basis for determining correctness. Inconsistent or ambiguous references could affect evaluation outcomes. Similarly, this study primarily uses binary verdicts which might overlook detailed aspects of responses that could be captured through more comprehensive evaluation criteria. Furthermore, while we conducted an error analysis of LLM judges and automatic metrics, there may be error cases that were not identified during our manual review, leaving gaps in understanding the full spectrum of evaluation inaccuracies. Finally, our study focuses exclusively on English, and the applicability of our approach to other languages, particularly morphologically rich or resource-scarce ones, remains unexplored.

## References

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2024. [Evaluating correctness and faithfulness of instruction-following models for question answering](#). *Transactions of the Association for Computational Linguistics*, 12:681–699.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2024. [Chateval: Towards better LLM-based evaluators through multi-agent debate](#). In *The Twelfth International Conference on Learning Representations*.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. [A survey on evaluation of large language models](#). *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. [Evaluating large language models trained on code](#).
- Yihan Chen, Benfeng Xu, Quan Wang, Yi Liu, and Zhendong Mao. 2024. [Benchmarking large language models on controllable generation under diversified instructions](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17808–17816.
- Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large language models be an alternative to human evaluations?](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.
- Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and Yiqun Liu. 2024. [Pre: A peer review based large language model evaluator](#).
- Domenic V Cicchetti and Alvan R Feinstein. 1990. High agreement but low kappa: Ii. resolving the paradoxes. *Journal of clinical epidemiology*, 43(6):551–558.
- Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. [All that’s ‘human’ is not gold: Evaluating human evaluation of generated text](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7282–7296, Online. Association for Computational Linguistics.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui

749	Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li,	Critic: Large language models can self-correct with	809
750	Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang	tool-interactive critiquing.	810
751	Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L.		
752	Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai	Rishav Hada, Varun Gumma, Adrian de Wynter,	811
753	Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai	Harshita Diddee, Mohamed Ahmed, Monojit Choud-	812
754	Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong	hury, Kalika Bali, and Sunayana Sitaram. 2024. <a href="#">Are</a>	813
755	Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan	<a href="#">large language model-based evaluators the solution</a>	814
756	Zhang, Minghua Zhang, Minghui Tang, Meng Li,	<a href="#">to scaling up multilingual evaluation?</a>	815
757	Miaojun Wang, Mingming Li, Ning Tian, Panpan		
758	Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen,	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and	816
759	Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan,	Weizhu Chen. 2021. <a href="#">Deberta: Decoding-enhanced</a>	817
760	Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen,	<a href="#">bert with disentangled attention</a> . In <i>International</i>	818
761	Shanghao Lu, Shangyan Zhou, Shanhuang Chen,	<i>Conference on Learning Representations</i> .	819
762	Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng		
763	Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	820
764	Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun,	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	821
765	T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu,	2021. <a href="#">Measuring massive multitask language under-</a>	822
766	Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao	<a href="#">standing</a> .	823
767	Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan		
768	Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin	Xinyu Hu, Mingqi Gao, Sen Hu, Yang Zhang, Yicheng	824
769	Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li,	Chen, Teng Xu, and Xiaojun Wan. 2024. <a href="#">Are llm-</a>	825
770	Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin,	<a href="#">based evaluators confusing nlg quality criteria?</a>	826
771	Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxi-		
772	ang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang,	Hui Huang, Yingqi Qu, Xingyuan Bu, Hongli Zhou,	827
773	Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang	Jing Liu, Muyun Yang, Bing Xu, and Tiejun Zhao.	828
774	Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng	2024a. <a href="#">An empirical study of llm-as-a-judge for llm</a>	829
775	Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,	<a href="#">evaluation: Fine-tuned judge model is not a general</a>	830
776	Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang,	<a href="#">substitute for gpt-4</a> .	831
777	Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,		
778	Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yu-	Hui Huang, Yingqi Qu, Hongli Zhou, Jing Liu, Muyun	832
779	jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,	Yang, Bing Xu, and Tiejun Zhao. 2024b. <a href="#">On the limi-</a>	833
780	Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu,	<a href="#">tations of fine-tuned judge models for llm evaluation</a> .	834
781	Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu,		
782	Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan,	Gautier Izacard and Edouard Grave. 2021. <a href="#">Leveraging</a>	835
783	Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean	<a href="#">passage retrieval with generative models for open do-</a>	836
784	Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao,	<a href="#">main question answering</a> . In <i>Proceedings of the 16th</i>	837
785	Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi-	<i>Conference of the European Chapter of the Associ-</i>	838
786	jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song,	<i>ation for Computational Linguistics: Main Volume</i> ,	839
787	Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu	pages 874–880, Online. Association for Computa-	840
788	Zhang, and Zhen Zhang. 2025. <a href="#">Deepseek-r1: Incen-</a>	tional Linguistics.	841
789	<a href="#">tivizing reasoning capability in llms via reinforce-</a>		
790	<a href="#">ment learning</a> .	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	842
		sch, Chris Bamford, Devendra Singh Chaplot, Diego	843
791	Ehsan Doostmohammadi, Oskar Holmström, and Marco	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	844
792	Kuhlmann. 2024. How reliable are automatic eval-	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	845
793	uation methods for instruction-tuned llms? <i>arXiv</i>	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	846
794	<i>preprint arXiv:2402.10770</i> .	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	847
		and William El Sayed. 2023. <a href="#">Mistral 7b</a> .	848
795	Owain Evans, Owen Cotton-Barratt, Lukas Finnve-		
796	den, Adam Bales, Avital Balwit, Peter Wills, Luca	Albert Q. Jiang, Alexandre Sablayrolles, Antoine	849
797	Righetti, and William Saunders. 2021. <a href="#">Truthful ai:</a>	Roux, Arthur Mensch, Blanche Savary, Chris	850
798	<a href="#">Developing and governing ai that does not lie</a> .	Bamford, Devendra Singh Chaplot, Diego de las	851
		Casas, Emma Bou Hanna, Florian Bressand, Gi-	852
799	Joseph L Fleiss and Jacob Cohen. 1973. The equiva-	anna Lengyel, Guillaume Bour, Guillaume Lam-	853
800	lence of weighted kappa and the intraclass correlation	ple, L��lio Renard Lavaud, Lucile Saulnier, Marie-	854
801	coefficient as measures of reliability. <i>Educational</i>	Anne Lachaux, Pierre Stock, Sandeep Subramanian,	855
802	<i>and psychological measurement</i> , 33(3):613–619.	Sophia Yang, Szymon Antoniak, Teven Le Scao,	856
803	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen,	Th��ophile Gervet, Thibaut Lavril, Thomas Wang,	857
804	Yujiu Yang, Nan Duan, and Weizhu Chen. 2024a.	Timoth��e Lacroix, and William El Sayed. 2024. <a href="#">Mix-</a>	858
805	<a href="#">Critic: Large language models can self-correct with</a>	<a href="#">tral of experts</a> .	859
806	<a href="#">tool-interactive critiquing</a> .		
807	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen,	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke	860
808	Yujiu Yang, Nan Duan, and Weizhu Chen. 2024b.	Zettlemoyer. 2017. <a href="#">Triviaqa: A large scale distantly</a>	861
		<a href="#">supervised challenge dataset for reading comprehen-</a>	862
		<a href="#">sion</a> .	863

864	Jaehun Jung, Faeze Brahman, and Yejin Choi. 2024.	Retrieval-augmented generation for knowledge-	922
865	<a href="#">Trust or escalate: Llm judges with provable guaran-</a>	intensive nlp tasks. In <i>Proceedings of the 34th Inter-</i>	923
866	<a href="#">tees for human agreement</a> .	<i>national Conference on Neural Information Process-</i>	924
867	Ehsan Kamalloo, Nouha Dziri, Charles Clarke, and	<i>ing Systems, NIPS '20</i> , Red Hook, NY, USA. Curran	925
868	Davood Rafiei. 2023. <a href="#">Evaluating open-domain ques-</a>	Associates Inc.	926
869	<a href="#">tion answering in the era of large language models</a> .	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad	927
870	In <i>Proceedings of the 61st Annual Meeting of the</i>	Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-	928
871	<i>Association for Computational Linguistics (Volume</i>	tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu,	929
872	<i>1: Long Papers)</i> , pages 5591–5606, Toronto, Canada.	Kai Shu, Lu Cheng, and Huan Liu. 2024a. <a href="#">From gen-</a>	930
873	Association for Computational Linguistics.	<a href="#">eration to judgment: Opportunities and challenges of</a>	931
874	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B.	<a href="#">llm-as-a-judge</a> .	932
875	Brown, Benjamin Chess, Rewon Child, Scott Gray,	Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia	933
876	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024b.	934
877	<a href="#">Scaling laws for neural language models</a> .	<a href="#">Llms-as-judges: A comprehensive survey on llm-</a>	935
878	Zachary Kenton, Noah Y. Siegel, János Kramár,	<a href="#">based evaluation methods</a> .	936
879	Jonah Brown-Cohen, Samuel Albanie, Jannis Bu-	Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan,	937
880	lian, Rishabh Agarwal, David Lindner, Yunhao Tang,	Hai Zhao, and Pengfei Liu. 2023. <a href="#">Generative judge</a>	938
881	Noah D. Goodman, and Rohin Shah. 2024. <a href="#">On scal-</a>	<a href="#">for evaluating alignment</a> .	939
882	<a href="#">able oversight with weak llms judging strong llms</a> .	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	940
883	Akbir Khan, John Hughes, Dan Valentine, Laura	Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and	941
884	Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward	Zhaopeng Tu. 2024. <a href="#">Encouraging divergent thinking</a>	942
885	Grefenstette, Samuel R. Bowman, Tim Rocktäschel,	<a href="#">in large language models through multi-agent debate</a> .	943
886	and Ethan Perez. 2024. <a href="#">Debating with more persua-</a>	In <i>Proceedings of the 2024 Conference on Empiri-</i>	944
887	<a href="#">sive llms leads to more truthful answers</a> .	<i>cal Methods in Natural Language Processing</i> , pages	945
888	Seungone Kim, Juyoung Suk, Shayne Longpre,	17889–17904, Miami, Florida, USA. Association for	946
889	Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham	Computational Linguistics.	947
890	Neubig, Moontae Lee, Kyungjae Lee, and Minjoon	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for auto-</a>	948
891	Seo. 2024. <a href="#">Prometheus 2: An open source language</a>	<a href="#">matic evaluation of summaries</a> . In <i>Text Summariza-</i>	949
892	<a href="#">model specialized in evaluating other language mod-</a>	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	950
893	<a href="#">els</a> .	Association for Computational Linguistics.	951
894	Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang,	952
895	Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and Xiaohang	Ruochen Xu, and Chenguang Zhu. 2023. <a href="#">G-eval:</a>	953
896	Dong. 2024. <a href="#">Better zero-shot reasoning with</a>	<a href="#">NLG evaluation using gpt-4 with better human align-</a>	954
897	<a href="#">role-play prompting</a> .	<a href="#">ment</a> . In <i>Proceedings of the 2023 Conference on</i>	955
898	Tom Kwiakowski, Jennimaria Palomaki, Olivia Red-	<i>Empirical Methods in Natural Language Processing</i> ,	956
899	field, Michael Collins, Ankur Parikh, Chris Alberti,	pages 2511–2522, Singapore. Association for Com-	957
900	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	putational Linguistics.	958
901	ton Lee, Kristina Toutanova, Llion Jones, Matthew	Yuxuan Liu, Tianchi Yang, Shaohan Huang, Zihan	959
902	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	Zhang, Haizhen Huang, Furu Wei, Weiwei Deng,	960
903	Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natu-</a>	Feng Sun, and Qi Zhang. 2024. <a href="#">Calibrating LLM-</a>	961
904	<a href="#">ral questions: A benchmark for question answering</a>	<a href="#">based evaluator</a> . In <i>Proceedings of the 2024 Joint</i>	962
905	<a href="#">research</a> . <i>Transactions of the Association for Compu-</i>	<i>International Conference on Computational Linguis-</i>	963
906	<i>tational Linguistics</i> , 7:452–466.	<i>tics, Language Resources and Evaluation (LREC-</i>	964
907	Md Tahmid Rahman Laskar, Sawsan Alqahtani, M Sai-	<i>COLING 2024)</i> , pages 2638–2656, Torino, Italia.	965
908	ful Bari, Mizanur Rahman, Mohammad Abdul-	ELRA and ICCL.	966
909	lah Matin Khan, Haidar Khan, Israt Jahan, Amran	Oscar Mañas, Benno Krojer, and Aishwarya Agrawal.	967
910	Bhuiyan, Chee Wei Tan, Md Rizwan Parvez, Enamul	2024. <a href="#">Improving automatic vqa evaluation using</a>	968
911	Hoque, Shafiq Joty, and Jimmy Huang. 2024. <a href="#">A sys-</a>	<a href="#">large language models</a> . In <i>Proceedings of the AAAI</i>	969
912	<a href="#">tematic survey and critical review on evaluating large</a>	<i>Conference on Artificial Intelligence</i> , volume 38,	970
913	<a href="#">language models: Challenges, limitations, and recom-</a>	pages 4171–4179.	971
914	<a href="#">mendations</a> . In <i>Proceedings of the 2024 Conference</i>	Mary L McHugh. 2012. Interrater reliability: the kappa	972
915	<i>on Empirical Methods in Natural Language Process-</i>	statistic. <i>Biochemia medica</i> , 22(3):276–282.	973
916	<i>ing</i> , pages 13785–13816, Miami, Florida, USA. As-	Meta AI. 2024. <a href="#">Introducing meta llama 3: The most</a>	974
917	sociation for Computational Linguistics.	<a href="#">capable openly available llm to date</a> . Meta AI Blog.	975
918	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	Accessed: 2024-07-25, 12:14:31 p.m.	976
919	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-		
920	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-		
921	täschel, Sebastian Riedel, and Douwe Kiela. 2020.		



977	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe,	Jake McNeil, David Medina, Aalok Mehta, Jacob	1039
978	Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-	Menick, Luke Metz, Andrey Mishchenko, Pamela	1040
979	moyer. 2022. <a href="#">Rethinking the role of demonstrations:</a>	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	1041
980	<a href="#">What makes in-context learning work?</a> In <i>Proceed-</i>	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	1042
981	<i>ings of the 2022 Conference on Empirical Methods in</i>	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	1043
982	<i>Natural Language Processing</i> , pages 11048–11064,	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	1044
983	Abu Dhabi, United Arab Emirates. Association for	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	1045
984	Computational Linguistics.	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	1046
985	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and	tista Parascandolo, Joel Parish, Emy Parparita, Alex	1047
986	Luke Zettlemoyer. 2020. <a href="#">AmbigQA: Answering am-</a>	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	1048
987	<a href="#">biguous open-domain questions</a> . In <i>Proceedings of</i>	man, Filipe de Avila Belbute Peres, Michael Petrov,	1049
988	<i>the 2020 Conference on Empirical Methods in Nat-</i>	Henrique Ponde de Oliveira Pinto, Michael, Poko-	1050
989	<i>ural Language Processing (EMNLP)</i> , pages 5783–	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	1051
990	5797, Online. Association for Computational Lin-	ell, Alethea Power, Boris Power, Elizabeth Proehl,	1052
991	guistics.	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	1053
992	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	Cameron Raymond, Francis Real, Kendra Rimbach,	1054
993	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	1055
994	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	1056
995	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	Girish Sastry, Heather Schmidt, David Schnurr, John	1057
996	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	Schulman, Daniel Selsam, Kyla Sheppard, Toki	1058
997	ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	1059
998	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	1060
999	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	1061
1000	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	1062
1001	man, Tim Brooks, Miles Brundage, Kevin Button,	lipo Petroski Such, Natalie Summers, Ilya Sutskever,	1063
1002	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	1064
1003	Carey, Chelsea Carlson, Rory Carmichael, Brooke	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	1065
1004	Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	1066
1005	Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	lipo Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	1067
1006	Chess, Chester Cho, Casey Chu, Hyung Won Chung,	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	1068
1007	Dave Cummings, Jeremiah Currier, Yunxing Dai,	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	1069
1008	Cory Decareaux, Thomas Degry, Noah Deutsch,	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	1070
1009	Damien Deville, Arka Dhar, David Dohan, Steve	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	1071
1010	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	Clemens Winter, Samuel Wolrich, Hannah Wong,	1072
1011	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	1073
1012	Simón Posada Fishman, Juston Forte, Isabella Ful-	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	1074
1013	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	1075
1014	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	1076
1015	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	Zheng, Juntang Zhuang, William Zhuk, and Barret	1077
1016	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	Zoph. 2023. <a href="#">Gpt-4 technical report</a> .	1078
1017	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	1079
1018	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	1080
1019	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	1081
1020	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	1082
1021	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	1083
1022	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-	1084
1023	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	1085
1024	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	1086
1025	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	1087
1026	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	man, Tim Brooks, Miles Brundage, Kevin Button,	1088
1027	Christina Kim, Yongjik Kim, Jan Hendrik Kirch-	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	1089
1028	ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo,	Carey, Chelsea Carlson, Rory Carmichael, Brooke	1090
1029	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	1091
1030	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	1092
1031	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	Chess, Chester Cho, Casey Chu, Hyung Won Chung,	1093
1032	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	Dave Cummings, Jeremiah Currier, Yunxing Dai,	1094
1033	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	Cory Decareaux, Thomas Degry, Noah Deutsch,	1095
1034	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	Damien Deville, Arka Dhar, David Dohan, Steve	1096
1035	Anna Makanju, Kim Malfacini, Sam Manning, Todor	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	1097
1036	Markov, Yaniv Markovski, Bianca Martin, Katie	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	1098
1037	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	Simón Posada Fishman, Juston Forte, Isabella Ful-	1099
1038	McKinney, Christine McLeavey, Paul McMillan,	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	1100
		Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	1101



1102	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	Zoph. 2024. <a href="#">Gpt-4 technical report</a> .	1165
1103	Gray, Ryan Greene, Joshua Gross, Shixiang Shane		
1104	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	1166
1105	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evalu-</a>	1167
1106	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	<a href="#">ation of machine translation</a> . In <i>Proceedings of the</i>	1168
1107	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	<i>40th Annual Meeting of the Association for Compu-</i>	1169
1108	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	<i>tational Linguistics</i> , pages 311–318, Philadelphia,	1170
1109	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	Pennsylvania, USA. Association for Computational	1171
1110	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	Linguistics.	1172
1111	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-		
1112	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	Keita Saito, Akifumi Wachi, Koki Wataoka, and Youhei	1173
1113	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	Akimoto. 2023. <a href="#">Verbosity bias in preference labeling</a>	1174
1114	Christina Kim, Yongjik Kim, Jan Hendrik Kirch-	<a href="#">by large language models</a> .	1175
1115	ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo,		
1116	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	Chenglei Si, Chen Zhao, and Jordan Boyd-Graber. 2021.	1176
1117	stantinidis, Kyle Kosic, Gretchen Krueger, Vishal	<a href="#">What’s in a name? answer equivalence for open-</a>	1177
1118	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	<a href="#">domain question answering</a> . In <i>Proceedings of the</i>	1178
1119	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	<i>2021 Conference on Empirical Methods in Natural</i>	1179
1120	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	<i>Language Processing</i> , pages 9623–9629, Online and	1180
1121	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	Punta Cana, Dominican Republic. Association for	1181
1122	Anna Makanju, Kim Malfacini, Sam Manning, Todor	Computational Linguistics.	1182
1123	Markov, Yaniv Markovski, Bianca Martin, Katie		
1124	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	Guangzhi Sun, Anmol Kagrecha, Potsawee Manakul,	1183
1125	McKinney, Christine McLeavey, Paul McMillan,	Phil Woodland, and Mark Gales. 2024. <a href="#">Skillaggre-</a>	1184
1126	Jake McNeil, David Medina, Aalok Mehta, Jacob	<a href="#">gation: Reference-free llm-dependent aggregation</a> .	1185
1127	Menick, Luke Metz, Andrey Mishchenko, Pamela	<i>arXiv preprint arXiv:2410.10215</i> .	1186
1128	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel		
1129	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	Tianxiang Sun, Junliang He, Xipeng Qiu, and Xuan-	1187
1130	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	jing Huang. 2022. <a href="#">BERTScore is unfair: On social</a>	1188
1131	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	<a href="#">bias in language model-based metrics for text ge-</a>	1189
1132	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	<a href="#">neration</a> . In <i>Proceedings of the 2022 Conference on</i>	1190
1133	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	<i>Empirical Methods in Natural Language Processing</i> ,	1191
1134	tista Parascandolo, Joel Parish, Emy Parparita, Alex	pages 3726–3739, Abu Dhabi, United Arab Emirates.	1192
1135	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	Association for Computational Linguistics.	1193
1136	man, Filipe de Avila Belbute Peres, Michael Petrov,		
1137	Henrique Ponde de Oliveira Pinto, Michael, Poko-	Aman Singh Thakur, Kartik Choudhary, Venkat Srinik	1194
1138	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	Ramayapally, Sankaran Vaidyanathan, and Dieuwke	1195
1139	ell, Alethea Power, Boris Power, Elizabeth Proehl,	Hupkes. 2024. <a href="#">Judging the judges: Evaluating align-</a>	1196
1140	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	<a href="#">ment and vulnerabilities in llms-as-judges</a> .	1197
1141	Cameron Raymond, Francis Real, Kendra Rimbach,		
1142	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yix-	1198
1143	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	uan Su, Aleksandra Piktus, Arkady Arkhangorodsky,	1199
1144	Girish Sastry, Heather Schmidt, David Schnurr, John	Minjie Xu, Naomi White, and Patrick Lewis. 2024.	1200
1145	Schulman, Daniel Selsam, Kyla Sheppard, Toki	<a href="#">Replacing judges with juries: Evaluating llm genera-</a>	1201
1146	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	<a href="#">tions with a panel of diverse models</a> .	1202
1147	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,		
1148	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry	1203
1149	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny	1204
1150	lippe Petroski Such, Natalie Summers, Ilya Sutskever,	Zhou, Quoc Le, and Thang Luong. 2023. <a href="#">Freshllms:</a>	1205
1151	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	<a href="#">Refreshing large language models with search engine</a>	1206
1152	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	<a href="#">augmentation</a> .	1207
1153	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-		
1154	lippe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	Tu Vu, Kalpesh Krishna, Salaheddin Alzubi, Chris Tar,	1208
1155	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	Manaal Faruqui, and Yun-Hsuan Sung. 2024. <a href="#">Foun-</a>	1209
1156	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	<a href="#">dational autoraters: Taming large language models</a>	1210
1157	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	<a href="#">for better automatic evaluation</a> .	1211
1158	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,		
1159	Clemens Winter, Samuel Wolrich, Hannah Wong,	Cunxiang Wang, Sirui Cheng, Qipeng Guo, Yuanhao	1212
1160	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	Yue, Bowen Ding, Zhikun Xu, Yidong Wang, Xi-	1213
1161	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	angkun Hu, Zheng Zhang, and Yue Zhang. 2024a.	1214
1162	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	<a href="#">Evaluating open-qa evaluation</a> . In <i>Proceedings of the</i>	1215
1163	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	<i>37th International Conference on Neural Information</i>	1216
1164	Zheng, Juntang Zhuang, William Zhuk, and Barret	<i>Processing Systems, NIPS ’23</i> , Red Hook, NY, USA.	1217
		Curran Associates Inc.	1218

1219	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu,	Lianghui Zhu, Xinggang Wang, and Xinlong Wang.	1274
1220	Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and	2023. <a href="#">Judgelm: Fine-tuned large language</a>	1275
1221	Zhifang Sui. 2023. <a href="#">Large language models are not</a>	<a href="#">models are scalable judges.</a> <i>arXiv preprint</i>	1276
1222	<a href="#">fair evaluators.</a>	<i>arXiv:2310.17631.</i>	1277
1223	Yuqi Wang, Lyuhao Chen, Songcheng Cai, Zhijian Xu,	<b>A Free-form Question-Answering</b>	1278
1224	and Yilun Zhao. 2024b. <a href="#">Revisiting automated evalu-</a>	In our experiments, we include AmbigQA (Min	1279
1225	<a href="#">ation for long-form table question answering.</a> In <i>Pro-</i>	et al., 2020), FreshQA (Vu et al., 2023), Hot-	1280
1226	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	potQA (Yang et al., 2018), Natural Ques-	1281
1227	<i>ods in Natural Language Processing</i> , pages 14696–	tions (Kwiatkowski et al., 2019), and Trivi-	1282
1228	14706, Miami, Florida, USA. Association for Com-	aQA (Joshi et al., 2017).	1283
1229	putational Linguistics.		
1230	Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao	• <b>AmbigQA:</b> Focuses on 14K ambiguous ques-	1284
1231	Song, Markus Freitag, William Wang, and Lei Li.	tions derived from NQ, requiring systems to	1285
1232	2023. <a href="#">INSTRUCTSCORE: Towards explainable text</a>	identify multiple valid interpretations and gen-	1286
1233	<a href="#">generation evaluation with automatic feedback.</a> In	erate disambiguated questions alongside cor-	1287
1234	<i>Proceedings of the 2023 Conference on Empirical</i>	responding answers.	1288
1235	<i>Methods in Natural Language Processing</i> , pages		
1236	5967–5994, Singapore. Association for Computa-	• <b>FreshQA:</b> A QA benchmark containing 600	1289
1237	tional Linguistics.	questions that consist of a diverse range of	1290
1238	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-	types, including those requiring fast-changing	1291
1239	gio, William W. Cohen, Ruslan Salakhutdinov, and	world knowledge and questions with false	1292
1240	Christopher D. Manning. 2018. <a href="#">Hotpotqa: A dataset</a>	premises that need debunking. It is regularly	1293
1241	<a href="#">for diverse, explainable multi-hop question answer-</a>	updated to reflect current information and is	1294
1242	<a href="#">ing.</a>	designed to evaluate the factual accuracy of	1295
1243	Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen,	LLMs in handling up-to-date and evolving	1296
1244	Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer,	knowledge.	1297
1245	Chao Huang, Pin-Yu Chen, Nitesh V Chawla, and	• <b>HotpotQA:</b> Contains 113K questions based	1298
1246	Xiangliang Zhang. 2024. <a href="#">Justice or prejudice? quan-</a>	on Wikipedia. It is designed to test multi-	1299
1247	<a href="#">tifying biases in llm-as-a-judge.</a>	hop reasoning, requiring connections across	1300
1248	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.	multiple paragraphs, and includes annotated	1301
1249	Weinberger, and Yoav Artzi. 2020. <a href="#">Bertscore: Evalu-</a>	supporting facts for evaluation.	1302
1250	<a href="#">ating text generation with BERT.</a> In <i>8th International</i>		
1251	<i>Conference on Learning Representations, ICLR 2020,</i>	• <b>Natural Questions (NQ):</b> Consists of real	1303
1252	<i>Addis Ababa, Ethiopia, April 26-30, 2020.</i> OpenRe-	user queries from Google Search, paired with	1304
1253	<i>view.net.</i>	Wikipedia articles. The dataset includes 307K	1305
1254	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	training examples annotated with both long	1306
1255	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	(paragraph) and short (entity-level) answers.	1307
1256	Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei	• <b>TriviaQA:</b> Features approximately 650K	1308
1257	Bi, Freda Shi, and Shuming Shi. 2023. <a href="#">Siren’s song</a>	trivia questions, with evidence sourced from	1309
1258	<a href="#">in the ai ocean: A survey on hallucination in large</a>	Wikipedia and web searches. These questions	1310
1259	<a href="#">language models.</a>	often require reasoning across multiple docu-	1311
1260	Yue Zhang, Ming Zhang, Haipeng Yuan, Shichun Liu,	ments for complex answer synthesis.	1312
1261	Yongyao Shi, Tao Gui, Qi Zhang, and Xuanjing		
1262	Huang. 2024. <a href="#">Llmeval: A preliminary study on</a>	We utilize the validation splits across multiple	1313
1263	<a href="#">how to evaluate large language models.</a> <i>Proceedings</i>	datasets: the standard validation split for Am-	1314
1264	<i>of the AAAI Conference on Artificial Intelligence,</i>	bigQA and Natural Questions, the “distractor” sub-	1315
1265	38(17):19615–19622.	set’s validation split for HotpotQA, and the “unfil-	1316
1266	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	tered.nocontext” subset’s validation split for Trivi-	1317
1267	Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin,	aQA. We randomly sampled 300 examples from	1318
1268	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	each dataset using Seed 42.	1319
1269	Joseph E. Gonzalez, and Ion Stoica. 2024. <a href="#">Judging</a>		
1270	<a href="#">llm-as-a-judge with mt-bench and chatbot arena.</a> In		
1271	<i>Proceedings of the 37th International Conference on</i>		
1272	<i>Neural Information Processing Systems, NIPS ’23,</i>		
1273	Red Hook, NY, USA. Curran Associates Inc.		

## B Human evaluation

This section provides detailed guidelines for human annotators responsible for evaluating the outputs of candidate LLMs. The goal is to ensure consistency and objectivity across all evaluations. These guidelines provide clear instructions for assessing each model’s response based on its alignment with the reference answer and contextual relevance.

### B.1 Guidelines

Dear Evaluator,

Thank you for your valuable contribution to this evaluation process. These guidelines outline the process for evaluating Large Language Model (LLM) outputs for the given tasks. As annotators, you will receive three components for each evaluation instance: the input question, reference answer(s), and the model’s response. Your task is to evaluate the responses independently and score them on a binary scale: ‘1’ for ‘True’ (correct) and ‘0’ for ‘False’ (incorrect).

A response warrants a score of ‘1’ when it demonstrates semantic equivalence with the reference answer, even if expressed through alternative phrasing or structure. This includes acceptable variations such as synonym usage and structural variations. Additional contextual information is acceptable as long as it doesn’t introduce errors.

Responses receive a score of ‘0’ when they contain factual errors, miss crucial elements from the reference answer, or demonstrate contextual misalignment. Partial answers that omit essential information should be marked incorrect, regardless of the accuracy of included content. When multiple reference answers are provided, a response is correct if it fully aligns with at least one reference. You are encouraged to use internet resources when needed to verify specific facts, terminology, or potential synonyms that may affect your evaluation decision. However, the reference answer should remain the primary basis for evaluation. Focus on whether the model’s response conveys the same core information as the reference answer. To maintain reliability, document any challenging cases requiring further discussion with other annotators.

### B.2 Inter human annotator agreement

We calculate Fleiss’ Kappa ( $\kappa$ ) (Fleiss and Cohen, 1973) to assess inter-rater reliability among human annotators. Table 4 and 5 show the inter-annotator agreement across models and tasks.

The results demonstrate high reliability, with Fleiss’ Kappa scores consistently above 0.93 for most tasks. The highest agreement is observed in Mixtral evaluations on HotpotQA ( $\kappa = 0.996$ ), and GPT on NQ-Open ( $\kappa = 0.990$ ). In FreshQA, which shows lower Kappa scores, the agreement among annotators remains high including 99.3% in GPT and 98.0% in Mixtral.

The percent agreement scores in Table 5 further confirm strong inter-annotator consistency. Most models achieve over 98% agreement across AmbigQA, HotpotQA, NQ-Open, and TriviaQA. However, DeepSeek exhibits lower agreement on NQ-Open (92.0%) and TriviaQA (90.0%). This indicates a variance in human ratings for these tasks.

## C Additional results

This section provides further results and analysis of conventional metrics and LLM-based evaluators. Table 6 illustrates the overall performance of candidate LLMs obtained through various evaluators. Unlike lexical matching and neural-based metrics, each LLM-as-a-judge indicates overall performance close to the human majority. Automatic metrics like EM severely underestimate the candidate LLMs’ performance. On the other hand, BERTScore tends to overestimate the performance.

EM underestimates performance because it requires a candidate’s response to exactly match one of the reference answers. This rigid, lexical approach fails to account for valid paraphrases, synonyms, or alternative expressions that convey the same meaning. In free-form QA tasks, where there can be multiple correct answers phrased in various ways, EM’s strict criteria often penalize responses that are semantically accurate but differ slightly in wording. As a result, it underestimates the true capabilities of candidate LLMs, leading to an incomplete assessment of their performance.

BERTScore relies on token-level semantic similarity, which rewards shallow lexical overlap rather than actual factual accuracy. For example, in cases where minor differences in wording (e.g., “The Treaty of Versailles was signed in 1919.” versus “The Treaty of Versailles ended in 1919.”) lead to opposing factual claims, BERTScore still scores the response high due to its emphasis on matching tokens (e.g., “signed” versus “ended”). Additionally, verbosity bias and threshold instability—where a default threshold (threshold = 0.5) is arbitrarily set—further inflate its raw accuracy. However,



LLMs	AmbigQA	FreshQA	HotpotQA	NQ-Open	TriviaQA
DeepSeek	0.975	0.949	0.986	0.889	0.456 ( $\kappa$ paradox)
Llama	0.945	0.962	0.973	0.985	0.935
GPT	0.989	0.973	0.982	0.990	0.948
Mixtral	0.981	0.945	0.996	0.977	0.936
Mistral	0.978	0.932	0.981	0.978	0.975

Table 4: Fleiss’ Kappa scores of human annotators across models and tasks.

LLMs	AmbigQA	FreshQA	HotpotQA	NQ-Open	TriviaQA
DeepSeek	99.0%	98.0%	99.7%	92.0%	90.0%
Llama	96.3%	98.0%	98.0%	99.0%	99.0%
GPT	99.3%	99.3%	98.7%	99.3%	99.0%
Mixtral	98.7%	98.0%	99.7%	98.3%	98.3%
Mistral	98.3%	97.0%	98.7%	98.3%	99.0%

Table 5: Human annotators percent agreement scores across candidate models and tasks.

when comparing raw accuracy with instance-level agreement metrics like Cohen’s kappa, which adjusts for class imbalance and penalizes asymmetric errors, the limitations of BERTScore become apparent.

### C.1 Impact of arbitration on dispute resolution

Figure 6 illustrates the impact of arbitration on resolving disagreements between primary judges. Arbitration, facilitated by GPT-3.5 as the tiebreaker, consistently improves performance across all tasks, particularly in FreshQA and TriviaQA, where Macro F1 increases by up to 21.5 points. In contrast, tasks like AmbigQA and HotpotQA, where primary judges initially exhibit stronger agreement, show smaller but still meaningful improvements. This highlights the critical role of arbitration in enhancing agreement and achieving closer alignment with ground truth, especially in cases of significant disagreement among primary judges.

Notably, evaluations of DeepSeek-v3 exhibit higher disagreement between Llama-3.1-70B and Mistral-7B, particularly in FreshQA (28.3%) and AmbigQA (25.7%). From our analysis, we did not find strong evidence explaining why DeepSeek-v3 leads to higher disagreement between the primary judges.

We observed substantial enhancements in Cohen’s Kappa scores across several tasks. For instance, as illustrated in Figure 7, in the AmbigQA Cohen’s Kappa increased from 0.881 to 0.911

for Llama. Similarly, in the same task, Cohen’s Kappa from 0.467 to 0.773 for candidate DeepSeek. These improvements demonstrate that the arbitration mechanism effectively enhances the reliability and consistency of evaluations, particularly in complex and ambiguous tasks where primary judges are more likely to disagree.

Some Cohen’s Kappa scores remain relatively low, particularly in FreshQA and DeepSeek-evaluated outputs. This is partially explained by the Kappa Paradox, where high agreement on extreme cases (e.g., clear correct/incorrect responses) and unbalanced class distributions can artificially lower the Kappa scores. In such cases, even when evaluators mostly agree, Cohen’s Kappa can appear lower than expected. Despite this, the arbitration process effectively mitigates inconsistencies, especially in tasks involving evolving knowledge and nuanced interpretations, such as FreshQA.

### C.2 Cost analysis

Human evaluation is the gold standard for assessing LLM-generated responses, but it is expensive and time-consuming. In our setup, we employed three human annotators who volunteered their efforts. However, if these annotators were compensated based on standard annotation rates, the cost of evaluating such outputs would be significantly higher. On the other hand, GPT-3.5-turbo, acting as an arbitrator in DAFE, incurs a cost that depends on the number of arbitration cases. In our evaluations, GPT-3.5-turbo was invoked 1,318 times, with



LLMs	Tasks	Evaluators							
		EM	BS	HM	DeepSeek	Llama	GPT	Mixtral	Mistral
DeepSeek	AmbigQA	56.3	80.0	84.3	86.3	73.7	75.0	62.3	93.3
	FreshQA	31.3	88.0	84.3	84.7	82.7	75.3	58.0	82.3
	HotpotQA	38.6	78.4	57.7	58.0	51.0	51.0	52.7	57.7
	NQ-Open	35.0	78.3	60.3	64.7	63.7	61.3	55.3	68.3
	TriviaQA	77.3	90.7	94.3	90.7	94.0	91.7	81.7	89.7
Llama	AmbigQA	42.3	63.0	67.0	64.0	65.3	64.7	63.0	66.0
	FreshQA	25.6	81.3	77.7	81.3	78.3	72.7	71.0	62.3
	HotpotQA	34.3	67.7	56.3	56.7	58.3	54.0	50.7	52.7
	NQ-Open	31.7	61.7	66.3	62.3	62.7	60.0	59.0	66.7
	TriviaQA	74.3	94.0	94.7	88.0	90.3	90.0	88.7	84.7
GPT	AmbigQA	49.7	78.0	71.7	70.3	70.0	68.0	65.7	71.0
	FreshQA	24.6	89.3	70.7	58.0	51.7	78.7	83.0	83.3
	HotpotQA	33.7	80.0	54.0	50.3	53.0	52.7	51.7	54.0
	NQ-Open	36.3	74.0	65.3	65.3	62.7	59.0	59.0	67.0
	TriviaQA	74.3	95.3	93.0	90.0	89.3	90.7	89.7	86.3
Mixtral	AmbigQA	37.7	70.3	61.7	58.7	57.3	62.0	59.3	61.7
	FreshQA	18.6	89.7	86.0	72.3	67.0	87.0	85.0	77.7
	HotpotQA	25.0	69.7	47.0	46.3	45.3	45.7	44.7	46.0
	NQ-Open	23.7	63.7	56.7	54.0	52.7	47.7	52.3	59.7
	TriviaQA	64.7	91.3	90.7	83.7	86.3	89.7	86.0	85.3
Mistral	AmbigQA	31.0	61.7	49.7	47.7	46.3	47.7	46.3	53.3
	FreshQA	15.6	80.0	81.7	60.7	59.0	83.7	84.0	86.0
	HotpotQA	23.7	64.7	40.0	39.3	39.0	38.0	37.0	39.0
	NQ-Open	22.7	60.0	46.0	41.3	40.0	43.3	41.3	50.0
	TriviaQA	62.0	94.3	83.7	78.0	81.3	81.0	79.7	85.0

Table 6: Raw performance of candidate LLMs across free-form QA tasks evaluated through various methods. HM represents Human Majority and BS denotes BERTScore.

an estimated total cost of \$0.59, which increases to \$5.40 if a 2048 max token setting is used (see Table 7). Since GPT-3.5 is only invoked when primary judges disagree, this selective arbitration substantially reduces overall evaluation expenses while maintaining high reliability in assessments. Rather than relying on a single model for evaluation, this multi-model arbitration approach enhances trust by mitigating biases and weaknesses inherent in any individual model.

By invoking the arbitrator only when disagreements occur (rather than evaluating all responses), DAFE reduces arbitration usage by 82–94% compared to a majority-voting system. This leads to:

- Over 90% fewer third-judge inferences, drastically lowering computational demand.
- Up to 95% cost savings by avoiding redundant model evaluations.
- Better scalability, making it practical for large-scale deployments

### C.3 DeepSeek as the arbitrator

To assess the impact of using DeepSeek as the arbitrator in DAFE, we conducted experiments by replacing GPT-3.5-turbo with DeepSeek. We evaluated this setup using different candidate models across multiple tasks. Specifically, we tested GPT-3.5 on TriviaQA, DeepSeek on NQ-Open, and Llama on FreshQA. The primary judges remained Llama and Mistral, and arbitration was invoked only in cases of disagreement. Our findings indicate that DeepSeek as the arbitrator achieves strong performance, with Macro-F1 scores of 91.23 on TriviaQA, 79.11 on NQ-Open, and 0.914 on FreshQA.

### C.4 Evaluating with one strong LLM-as-a-judge

While a single state-of-the-art evaluator can achieve strong performance in many cases, the dual-LLM framework remains critical for ensuring robustness, particularly in high-stakes or ambiguous scenarios.

To explore the potential of a more powerful single LLM, we evaluated GPT-3.5-turbo on HotpotQA and TriviaQA using GPT-4o as a judge.

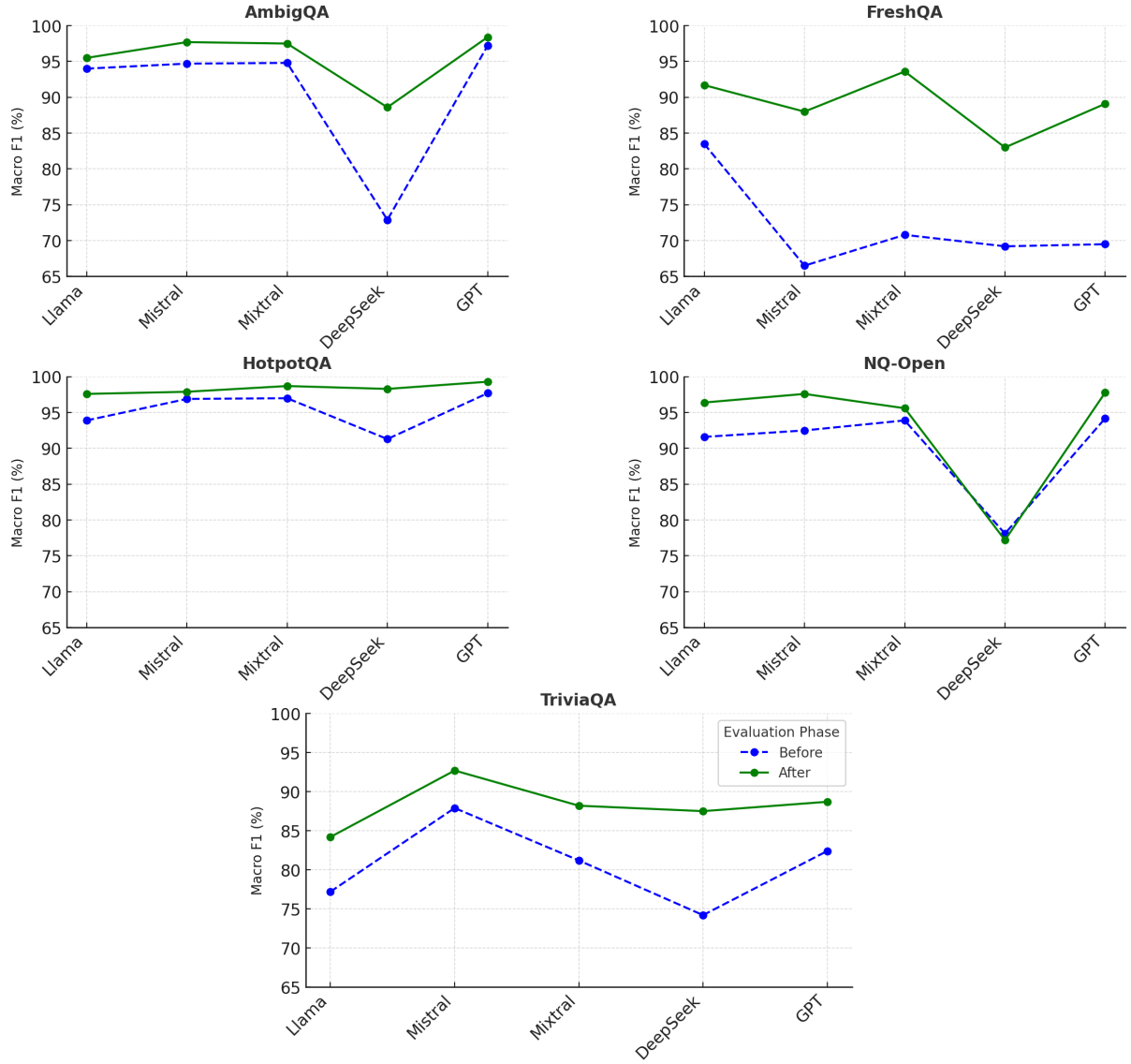


Figure 6: Impact of arbitration on disagreements between primary judges. Note that we used Llama-3.1-70B and Mistra 7B as primary judges. GPT-3.5-turbo is only utilized when disagreements are found. The models given in the figure are candidate LLMs which generate outputs for the given tasks and are then evaluated through DAFE.

With this configuration, GPT-4o as the evaluator achieved a Macro-F1 score of 0.946 on HotpotQA, demonstrating its exceptional capability. However, the same GPT-4o judge achieved only 0.784 on TriviaQA, which falls short of DAFE’s performance of 0.887. This shows that even the most advanced models show inconsistencies when evaluating free-form QA. This is particularly critical in precision-sensitive domains where minor errors can have outsized consequences.

In such settings, DAFE’s ensemble approach acts as a safeguard. When employing DAFE with GPT-3.5-turbo as the arbitrator, we achieved an even higher Macro-F1 of 0.984 on HotpotQA, surpassing the performance of a single GPT-4o. In-

terestingly, when we experimented with DeepSeek as the arbitrator in DAFE, performance remained strong at 0.963 Macro-F1, indicating that DAFE’s benefits are not solely tied to a specific arbitrator model.

### C.5 Majority voting-based evaluation

We conducted additional experiments utilizing a traditional majority voting approach for evaluating candidate LLM performance. In this setup, we employed three LLM judges of equal weight: Llama, GPT-3.5, and Mistral to evaluate candidate models generated response. For every evaluation instance, each judge provided an independent binary verdict (True or False). The final decision is determined

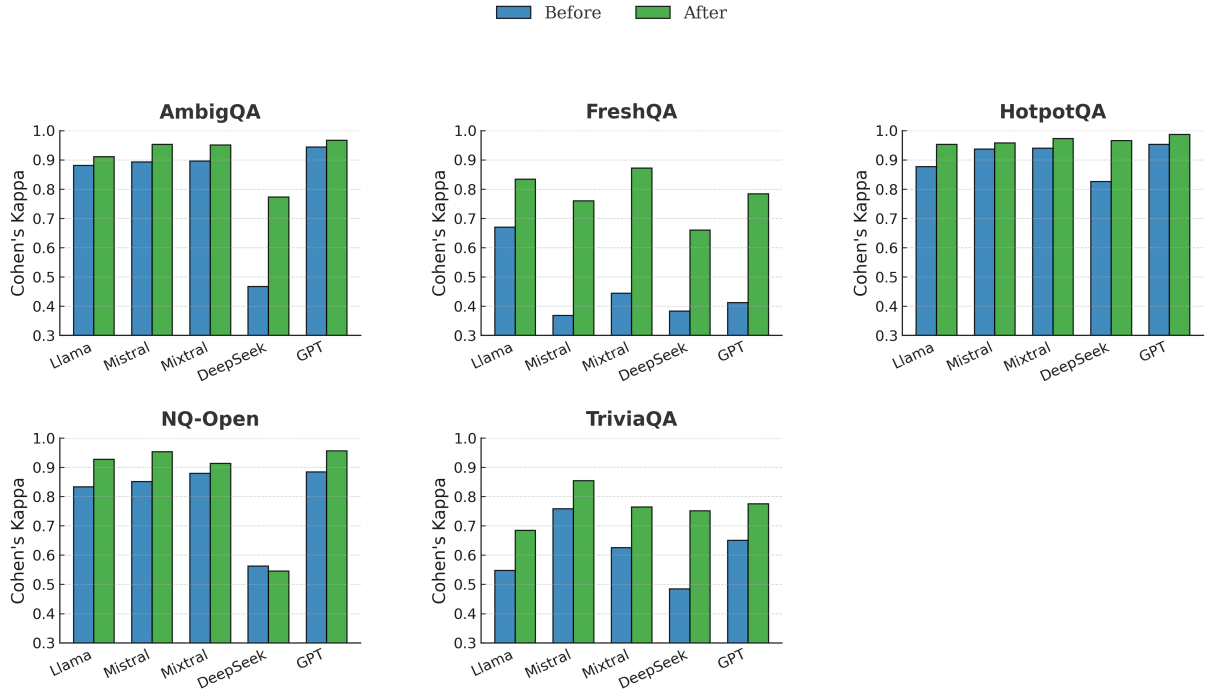


Figure 7: Comparison of Cohen’s kappa scores before and after arbitration (GPT-3.5-turbo as arbitrator). The performance is illustrated across candidate LLMs and tasks.

through a simple majority vote across these three verdicts.

As presented in Table 8, DAFE matches or closely approaches the Macro F1 and Cohen’s Kappa scores of the three-judge majority across almost all tasks and candidate LLMs. For example, on HotpotQA, evaluating candidate Llama with DAFE achieves a Macro F1 of 97.6% (compared to 97.6% for majority voting) and a Cohen’s Kappa of 0.95, while for GPT-3.5 on AmbigQA, DAFE reaches a Macro F1 of 98.4% (versus 98.3% for majority voting), indicating a negligible performance difference. Even in high-disagreement tasks like TriviaQA, where the primary judges (e.g., Mistral) disagree 20.3% of the time, DAFE retains strong alignment (with a Macro F1 of 92.7 compared to 93.5 for majority voting). Minor deviations, such as the one observed for candidate Mixtral on TriviaQA (DAFE’s Macro F1 = 0.88 vs. 0.95 for majority voting), reflect rare instances where both the primary judges and the arbitrator make errors, yet these outliers are substantially outweighed by the computational savings offered by selective arbitration.

## C.6 Impact of prompt variations

The effectiveness and consistency of LLM-based evaluation are significantly influenced by prompt

design. Variations in prompt structure, reasoning order, explanation requirements, and task-specific examples can lead to notable differences in model verdicts. To analyze the robustness of the LLM judges in free-form QA, we conducted ablation studies on different prompt variations using Mistral as the candidate model and GPT as the judge.

### C.6.1 Consistency in judgment across multiple trials

LLMs generate random text even at a temperature of 0. To assess whether this affects evaluation consistency, we repeated the same evaluation task five times for 100 Mistral-generated responses for HotpotQA.

- **Verdict stability:** GPT produced identical True/False verdicts in 100% of cases. This suggest that its binary decision-making process remains stable even across multiple trials.
- **Explanation variability:** While verdicts remained consistent, the rationales and explanations provided by GPT across trials, often cited different supporting facts for the same judgment.

### C.6.2 Few-shot vs. zero-shot prompting

We investigated the impact of few-shot prompting where we included three **task-specific examples**

Table 7: Cost-efficiency analysis of DAFE: Summary of disagreement rates and tiebreaker usage across candidate models and tasks

Candidate LLMs	Tasks	Samples	Disagreement Rates (%)	Tiebreaker Usage
DeepSeek	AmbigQA	300	25.7	77
	FreshQA	300	28.3	85
	HotpotQA	300	10.7	32
	NQ-Open	300	12.0	36
	TriviaQA	300	14.3	43
Llama	AmbigQA	300	10.0	30
	FreshQA	300	31.3	94
	HotpotQA	300	13.0	39
	NQ-Open	300	18.0	54
	TriviaQA	300	17.0	51
GPT	AmbigQA	300	7.0	21
	FreshQA	300	44.3	133
	HotpotQA	300	5.7	17
	NQ-Open	300	13.0	39
	TriviaQA	300	15.7	47
Mixtral	AmbigQA	300	9.0	27
	FreshQA	300	37.3	112
	HotpotQA	300	4.7	14
	NQ-Open	300	13.0	39
	TriviaQA	300	17.0	51
Mistral	AmbigQA	300	11.7	35
	FreshQA	300	39.7	119
	HotpotQA	300	6.0	18
	NQ-Open	300	14.7	44
	TriviaQA	300	20.3	61
<b>Total</b>		<b>7500</b>		<b>1318</b>

in the prompt to guide the judge’s decision-making process. We found that adding few-shot examples resulted in a 2% increase in Macro-F1 scores. However, few-shot prompting introduced rigid decision patterns—the model sometimes over-applied reasoning from the examples rather than adapting flexibly to novel cases. For instance, multi-hop reasoning cases from HotpotQA, the judge model consistently followed the structure of the provided examples, even when the correct reasoning required a different approach.

### C.6.3 Explanation requirement: Binary verdict vs. justification-based evaluation

To test whether requiring the model to generate explanations alongside verdicts improves judgment reliability, we compared two settings:

- **Binary verdict-only evaluation:** The model was instructed to provide only a True/False response without any explanation.
- **Justification-based evaluation:** The model was required to explain its reasoning before delivering the final verdict.

We found that:

- **Higher verdict volatility in verdict-only mode:** When explanations were removed, 13% of verdicts changed between repeated evaluations of the same responses.
- **Reduced alignment with human judgment:** Cohen’s Kappa agreement with human annotators dropped from 0.95 to 0.72, highlighting that rationale-based prompts lead to more stable and accurate decisions.



Table 8: Comparison between Majority Voting (Llama+GPT-3.5+Mistral) and DAFE (GPT-3.5 as arbitrator). For each candidate LLM and task, the table reports Macro F1 and Cohen’s Kappa scores under Majority Voting, the disagreement rate (in %), and the corresponding scores using DAFE.

Candidate LLM	Task	Majority Voting		Disagreement (%)	DAFE	
		Macro F1	Kappa		Macro F1	Kappa
Llama	AmbigQA	95.5	0.91	10.0	95.5	0.91
	HotpotQA	97.6	0.95	13.0	97.6	0.95
	NQ-Open	96.3	0.93	18.0	96.4	0.92
	TriviaQA	84.1	0.68	17.0	84.2	0.68
GPT	AmbigQA	98.3	0.97	7.0	98.4	0.96
	HotpotQA	99.3	0.99	5.7	99.3	0.98
	NQ-Open	97.8	0.96	13.0	97.8	0.95
	TriviaQA	90.5	0.81	15.7	88.7	0.77
Mistral	AmbigQA	98.9	0.98	9.0	97.5	0.95
	HotpotQA	98.6	0.97	4.7	98.7	0.97
	NQ-Open	98.3	0.97	13.0	95.6	0.91
	TriviaQA	95.0	0.90	17.0	88.2	0.76
Mistral	AmbigQA	97.6	0.95	11.7	97.7	0.95
	HotpotQA	97.9	0.96	6.0	97.9	0.95
	NQ-Open	97.6	0.95	14.7	97.6	0.95
	TriviaQA	93.5	0.87	20.3	92.7	0.85

#### C.6.4 Reason-first vs. verdict-first prompting

In the verdict-first approach, the model is instructed to provide a True/False answer before justifying its decision, whereas in the reason-first approach, the model is asked to generate reasoning first and then conclude with a verdict. Experimental results showed no significant difference in accuracy or agreement scores between these two formats.

#### C.7 G-Eval: reference-free evaluation of free-form question-answering

Existing LLM-based evaluators such as G-Eval (Liu et al., 2023) are designed for reference-free, subjective tasks (e.g., summarization, dialogue), where evaluation criteria (e.g., coherence, fluency) are inherently ambiguous and scored on Likert scales. These frameworks prioritize qualitative judgments rather than binary factual correctness. In contrast, DAFE is explicitly tailored for reference-dependent, objective evaluation in free-form QA, where answers are either factually correct or incorrect based on alignment with explicit ground-truth references.

To validate this distinction, we tailored G-eval based method to investigate the capability of LLM-as-a-judge in reference-free settings. In this setting,

we modify the evaluation prompt by excluding the reference answer  $r$  and directly prompted the evaluator model as  $P = \{x, \bar{y}\}$  along with instructions such as correctness.

The performance of LLM-as-a-judge drastically changes in reference-free settings. Without access to the ground truth references, we observe a stark decline in evaluation capability across all models (see Table 9 and 10 values in blue). This systematic deterioration spans all tasks and model combinations, though its severity varies by context. HotpotQA, with its demands for complex reasoning, exemplifies this challenge most clearly. The substantial gap between reference-based and reference-free evaluation underscores the crucial role of reference answers in reliable assessment.

#### C.8 DAFE in multi-reference answers

DAFE explicitly accommodates multiple gold reference answers by incorporating all available references into the judge LLM’s prompt during evaluation. For datasets like AmbigQA and TriviaQA, where questions often have multiple valid answers (e.g., synonyms, rephrased answers, or alternative factual representations), DAFE aggregates all reference answers into the judge’s input prompt (e.g.,

Candidate LLMs	Tasks	Evaluators						
		EM	BERTScore	Human Majority	Llama-3.1-70B	GPT-3.5-turbo	Mixtral-8x7B	Mistral-7B
Llama-3.1-70B	AmbigQA	42.3	63.0	67.0	65.3 [83.3]	64.7 [84.7]	63.0 [76.0]	66.0 [80.3]
	HotpotQA	34.3	67.7	56.3	58.3 [81.0]	54.0 [81.0]	50.7 [67.3]	52.7 [69.3]
	NQ-Open	31.7	61.7	66.3	62.7 [89.0]	60.0 [89.3]	59.0 [81.0]	66.7 [81.0]
	TriviaQA	74.3	94.0	94.7	90.3 [90.3]	90.0 [90.3]	88.7 [89.0]	84.7 [84.0]
GPT-3.5	AmbigQA	49.7	78.0	71.7	70.0 [79.0]	68.0 [81.0]	65.7 [79.0]	71.0 [84.3]
	HotpotQA	33.7	80.0	54.0	53.0 [85.3]	52.7 [85.7]	51.7 [82.3]	54.0 [86.3]
	NQ-Open	36.3	74.0	65.3	62.7 [83.7]	59.0 [90.7]	59.0 [87.0]	67.0 [89.7]
	TriviaQA	74.3	95.3	93.0	89.3 [89.0]	90.7 [88.7]	89.7 [90.3]	86.3 [84.3]
Mixtral-8x7B	AmbigQA	37.7	70.3	61.7	57.3 [74.7]	62.0 [82.3]	59.3 [79.7]	61.7 [80.7]
	HotpotQA	25.0	69.7	47.0	45.3 [80.0]	45.7 [84.7]	44.7 [72.0]	46.0 [78.0]
	NQ-Open	23.7	63.7	56.7	52.7 [81.7]	47.7 [90.3]	52.3 [85.7]	59.7 [89.7]
	TriviaQA	64.7	91.3	90.7	86.3 [85.7]	89.7 [89.0]	86.0 [86.7]	85.3 [86.0]
Mistral-7B	AmbigQA	31.0	61.7	49.7	46.3 [61.0]	47.7 [78.7]	46.3 [74.7]	53.3 [85.0]
	HotpotQA	23.7	64.7	40.0	39.0 [64.3]	38.0 [83.3]	37.0 [62.0]	39.0 [77.0]
	NQ-Open	22.7	60.0	46.0	40.0 [72.3]	43.3 [85.7]	41.3 [78.0]	50.0 [92.3]
	TriviaQA	62.0	94.3	83.7	81.3 [80.7]	81.0 [81.0]	79.7 [80.7]	85.0 [84.7]

Table 9: Overall performance (Accuracy) of candidate LLMs across free-form QA tasks. Values [in blue] represent LLM-as-a-judge in the reference-free mood.

Candidate LLMs	Tasks	Evaluators					
		EM	BERTScore	Llama-3.1-70B	GPT-3.5-turbo	Mixtral-8x7B	Mistral-7B
Llama-3.1-70B	AmbigQA	0.744	0.641	0.944 [0.629]	0.922 [0.604]	0.912 [0.669]	0.929 [0.631]
	HotpotQA	0.778	0.745	0.939 [0.628]	0.949 [0.574]	0.910 [0.665]	0.916 [0.640]
	NQ-Open	0.653	0.718	0.916 [0.606]	0.896 [0.560]	0.907 [0.639]	0.869 [0.622]
	TriviaQA	0.612	0.782	0.772 [0.772]	0.717 [0.628]	0.695 [0.678]	0.640 [0.633]
GPT-3.5	AmbigQA	0.792	0.622	0.972 [0.686]	0.949 [0.603]	0.930 [0.596]	0.927 [0.553]
	HotpotQA	0.794	0.623	0.977 [0.566]	0.987 [0.521]	0.936 [0.543]	0.966 [0.494]
	NQ-Open	0.703	0.606	0.942 [0.671]	0.911 [0.544]	0.911 [0.601]	0.914 [0.536]
	TriviaQA	0.646	0.681	0.824 [0.817]	0.700 [0.690]	0.789 [0.760]	0.730 [0.701]
Mixtral-8x7B	AmbigQA	0.760	0.666	0.948 [0.704]	0.891 [0.636]	0.955 [0.654]	0.944 [0.622]
	HotpotQA	0.761	0.657	0.970 [0.587]	0.966 [0.470]	0.930 [0.582]	0.970 [0.577]
	NQ-Open	0.650	0.649	0.939 [0.652]	0.863 [0.517]	0.950 [0.590]	0.908 [0.529]
	TriviaQA	0.625	0.695	0.812 [0.800]	0.803 [0.754]	0.838 [0.818]	0.716 [0.725]
Mistral-7B	AmbigQA	0.792	0.622	0.947 [0.730]	0.947 [0.627]	0.947 [0.628]	0.930 [0.523]
	HotpotQA	0.796	0.673	0.969 [0.649]	0.951 [0.478]	0.947 [0.680]	0.969 [0.578]
	NQ-Open	0.726	0.639	0.925 [0.652]	0.919 [0.515]	0.939 [0.597]	0.920 [0.433]
	TriviaQA	0.718	0.608	0.879 [0.881]	0.863 [0.840]	0.822 [0.846]	0.735 [0.744]

Table 10: Performance (Macro F1) of various evaluators across candidate LLMs and tasks. Values [in blue] represent LLM-as-a-judge in the reference-free mode.

concatenating them as a comma-separated list).

This design ensures that the judge evaluates the candidate’s output against the full spectrum of acceptable answers, mirroring the human evaluation protocol, where annotators are instructed to mark a response as correct if it aligns with any reference answer. However, as presented in our paper, LLM-based judges encounter challenges with multiple reference answers. This confusion is particularly evident in TriviaQA, where multiple reference answers introduce difficulties for the judges to recognize and evaluate a range of correct responses.

### C.9 Analysis of automatic metrics

Figures 8, 9, 10, and 11 illustrate the fundamental trade-offs in automatic metrics. In TriviaQA, where multiple normalized reference answers exist, EM achieves impressive true positives (61.7-74.3%)

compared to HotpotQA (23.0-34.3%) which contains single reference answers. EM’s near-zero false positives across tasks (0-0.7%) stem from its strict string matching – it only flags matches when answers are identical to references. Our error analysis found three primary causes of such rare false positives including preprocessing errors, where character normalization removes crucial distinctions, and reference ambiguities, where incomplete or ambiguous references lead to incorrect matches. Additionally, a semantic mismatch occurs when the EM incorrectly labels a prediction as true by matching text without considering its context. For instance, despite their different contextual meanings, EM wrongly marks a match between a model prediction of “1944” (describing the start of a war) and a reference answer containing “1944” (representing the end of the war).

EM string-matching guarantees high precision and makes EM particularly effective when exact wording is crucial, such as mathematical problems. However, its rigid criteria also result in substantial false negatives (17.0-34.7%). These false negatives primarily occur when the candidate LLM generates semantically correct responses that differ from references in format or expression. Common cases include synonym usage and paraphrases, structural variations in phrasing (e.g., “School of Medicine at Harvard” vs. “Harvard Medical School”), granularity discrepancies where answers differ in levels of detail from references (e.g., answering “British writer” instead of “William Shakespeare”), and partial matches that contain valid information but don’t exactly mirror the reference.

Unlike EM, BERTScore offers advantages in capturing semantic similarities. In TriviaQA, it gains high true positive rates (81.3-92.0%) with relatively low false positives (2.0-13.0%). BERTScore’s performance varies significantly across tasks and is influenced by its sensitivity to the threshold setting. In HotpotQA, where answers require multi-hop reasoning, true positives reach 36.0-50.3%, with an increase in false positives (17.7-29.7%). A similar pattern appears in NQ-Open, with true positives of 43.3-53.0% and false positives of 10.7-21.0%. Its tendency toward false positives indicates that relying solely on embedding similarity often accepts answers that are contextually related but factually incorrect. The false positives emerge through semantic drift (where similar embeddings yield false matches), contextual misalignment (where word meanings shift based on context), and threshold instability (where similarity cutoffs fail to distinguish subtle semantic differences). Additionally, false positives emerge due to the verbose responses where additional content artificially increases similarity scores.

## D Prompting

In our main experiment, we performed zero-shot prompting in the following two stages.

### D.1 Prompting Candidate LLMs

We prompted candidate LLMs (see Figure 12) to record generations for each task. We set the same role and prompt structure for each candidate model to ensure the reproducibility of our results. Figure 13 shows the candidate GPT-3.5-turbo response at zero temperature for the input given in Figure 12.

### D.2 Prompting LLM Judges

We prompted LLMs-as-judges to perform the evaluation (see Figure 14). In Figure 15, judge Llama-3.1-70B evaluating candidate GPT-3.5-turbo.

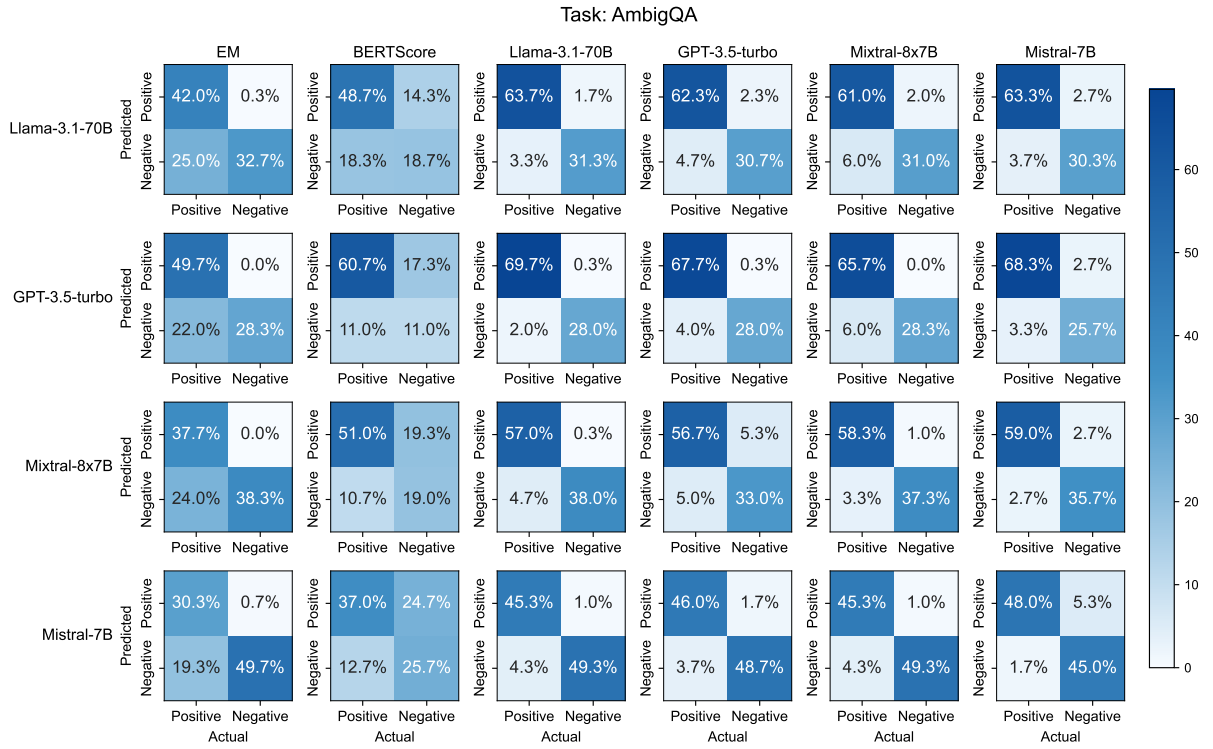


Figure 8: Confusion matrices comparing the performance of automatic metrics (EM, BERTScore) and individual LLM judges on AmbigQA.

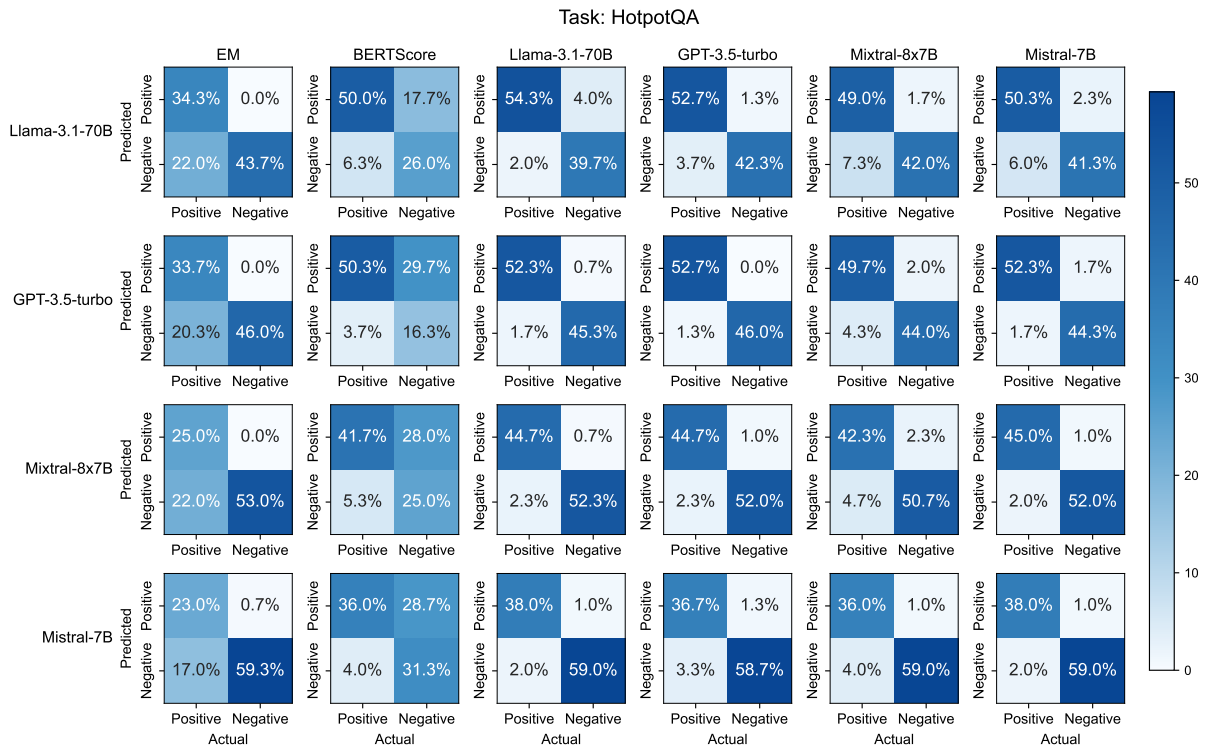


Figure 9: Confusion matrices comparing the performance of automatic metrics (EM, BERTScore) and individual LLM judges on HotpotQA.



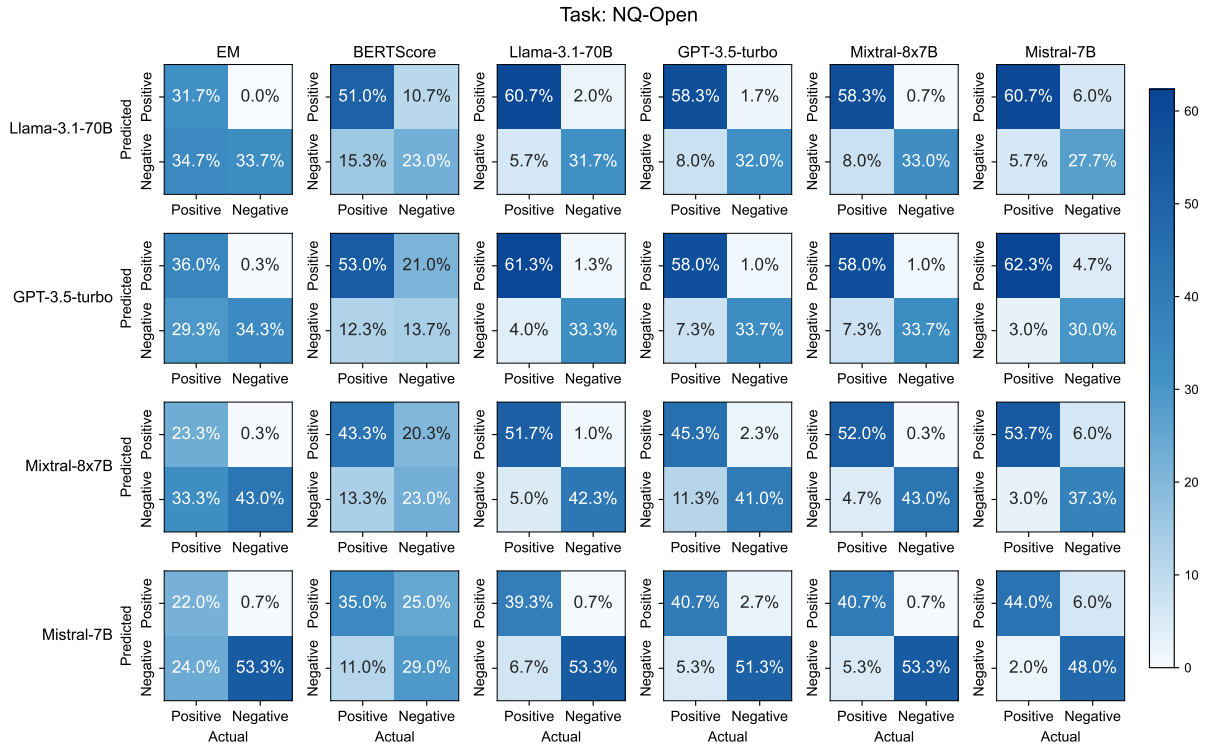


Figure 10: Confusion matrices comparing the performance of automatic metrics (EM, BERTScore) and individual LLM judges on NQ-Open.

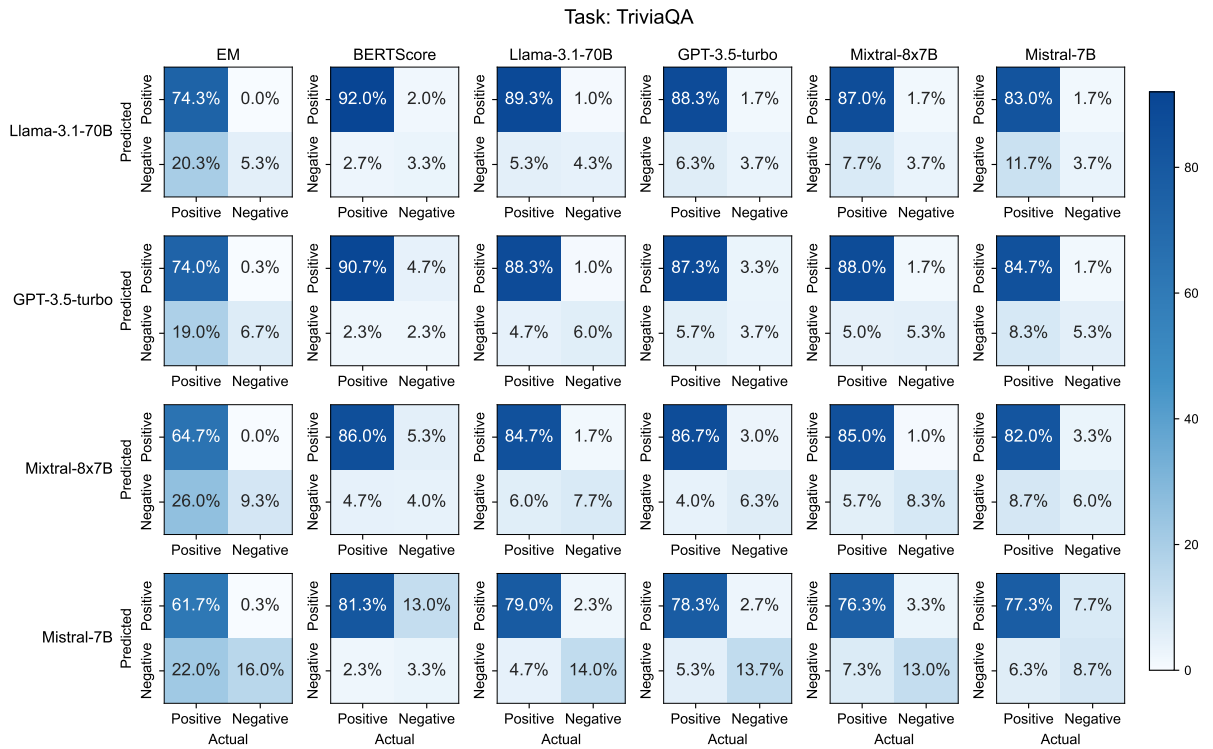


Figure 11: Confusion matrices comparing the performance of automatic metrics (EM, BERTScore) and individual LLM judges on TriviaQA.

You are a helpful assistant. Which comic book was also written by the writer of Crossed?

Figure 12: Prompting candidate GPT-3.5-turbo to elicit outputs for HotpotQA.

Crossed was written by Garth Ennis. The Boys is a comic book series also written by Garth Ennis. So the comic book that was also written by the writer of Crossed is The Boys.

Figure 13: Candidate GPT-3.5-turbo response to the input given in Figure 12.

You are a helpful assistant acting as an impartial judge. You will be given a Question and a Proposed Answer. Your task is to judge whether the Proposed Answer is correct by comparing it to the Reference Answer. If the Proposed Answer is correct, choose 'True', otherwise choose 'False'. Provide a brief explanation for your decision.

Question: Which comic book was also written by the writer of Crossed?

Provided Answer: Crossed was written by Garth Ennis. The Boys is a comic book series also written by Garth Ennis. So the comic book that was also written by the writer of Crossed is The Boys. So the answer is: The Boys.

Reference Answer: the Vertigo series "Preacher"

Evaluation:

Provide your response in the following format:

Decision: [True/False]

Explanation: [Your brief explanation]

Figure 14: Prompting judge Llama-3.1-70B for evaluating candidate GPT-3.5-turbo.

**Decision:** False

Explanation: The Proposed Answer incorrectly identifies "The Boys" as the comic book written by the writer of Crossed. While it is true that Garth Ennis wrote Crossed, the Proposed Answer fails to consider other works written by Garth Ennis. The Reference Answer, "Preacher", is indeed another comic book series written by Garth Ennis, specifically for the Vertigo imprint.

Figure 15: Llama-3.1-70B-Judge verdict on the candidate GPT-3.5-turbo output.