

A Systematic Survey and Critical Review on Evaluating Large Language Models: Challenges, Limitations, and Recommendations

Md Tahmid Rahman Laskar^{†,‡,*}, Sawsan Alqahtani[§], M Saiful Bari^{¶,*}
Mizanur Rahman^{†,*}, Mohammad Abdullah Matin Khan[‡], Haidar Khan[¶]
Israt Jahan[†], Md Amran Hossen Bhuiyan[†], Chee Wei Tan[‡], Md Rizwan Parvez[§]
Enamul Hoque[†], Shafiq Joty^{‡,°,*}, Jimmy Xiangji Huang^{†,*}

[†]York University, [§]Princess Nourah Bint Abdulrahman University, [‡]Nanyang Technological University,

[¶]National Center for AI, Saudi Arabia, [§]Qatar Computing Research Institute (QCRI),

[°]Dialpad Canada Inc., ^{*}Royal Bank of Canada, [°]Salesforce Research

Abstract

Large Language Models (LLMs) have recently gained significant attention due to their remarkable capabilities in performing diverse tasks across various domains. However, a thorough evaluation of these models is crucial before deploying them in real-world applications to ensure they produce reliable performance. Despite the well-established importance of evaluating LLMs in the community, the complexity of the evaluation process has led to varied evaluation setups, causing inconsistencies in findings and interpretations. To address this, we systematically review the primary challenges and limitations causing these inconsistencies and unreliable evaluations in various steps of LLM evaluation. Based on our critical review, we present our perspectives and recommendations to ensure LLM evaluations are reproducible, reliable, and robust.

1 Introduction

The evolution of LLMs has transitioned from simple generative models predicting the next word to advanced systems capable of following instructions and solving complex problems (Zhao et al., 2023a). Early models like GPT (Radford et al., 2018) could generate coherent text but were limited to simple tasks, whereas instruction-tuned LLMs (Chung et al., 2022; Ouyang et al., 2022) like ChatGPT¹ greatly enhanced their versatility and ability to execute specific commands. This shift has revolutionized the development of real-world applications powered by LLMs.

With the advancements and broad applicability of LLMs, it is essential to properly evaluate them to ensure they are safe to use. This is indeed important not only for academic benchmarks

but also for business use cases. Consequently, understanding the bottlenecks of current evaluation methods, and developing strategies to address these challenges are crucial for standardizing evaluations and enabling reliable use of LLMs in practical applications. Nonetheless, evaluating LLMs is as complex and resource-intensive as their development, involving multiple levels or aspects.

Existing reviews (Chang et al., 2024; Guo et al., 2023b; Liang et al., 2022; Minaee et al., 2024; Zhuang et al., 2023) related to the evaluation of LLMs often focus only on benchmark tasks, datasets, and evaluation criteria, neglecting the broader complexities. This oversight can undermine the reliability of evaluation by ignoring issues like robustness and reproducibility. While some recent studies (Balloccu et al., 2024; Mao et al., 2023) have investigated data contamination (Ravaut et al., 2024) and evaluation malpractices in LLM evaluation, their focus is limited to only assessing ChatGPT, overlooking other LLMs, as well as the entire evaluation pipeline.

More recently, Biderman et al. (2024) discussed the reproducibility problem in existing evaluations of LLMs and introduced a library to address this. However, their work lacks comprehensive discussions on how aspects like reliability or robustness impact LLM evaluation and how to address them. Hence, existing LLM evaluation studies often focus on individual aspects in a scattered manner, resulting in findings that are only sparsely useful.

To mitigate this gap, this paper brings together the discussions to address the fundamental challenges and limitations in LLM evaluations that emerge from diverse evaluation setups. First, we craft a schematic workflow of the evaluation pipeline in practical settings (presented in Section 2) for a systematic study. We then examine each step in the evaluation workflow, uncovering various inconsistencies and decision-making complexities affecting repro-

*Corresponding Emails: {tahmid20, jhuang}@yorku.ca, {bari0001, srjoty}@ntu.edu.sg

¹<https://openai.com/index/chatgpt/>

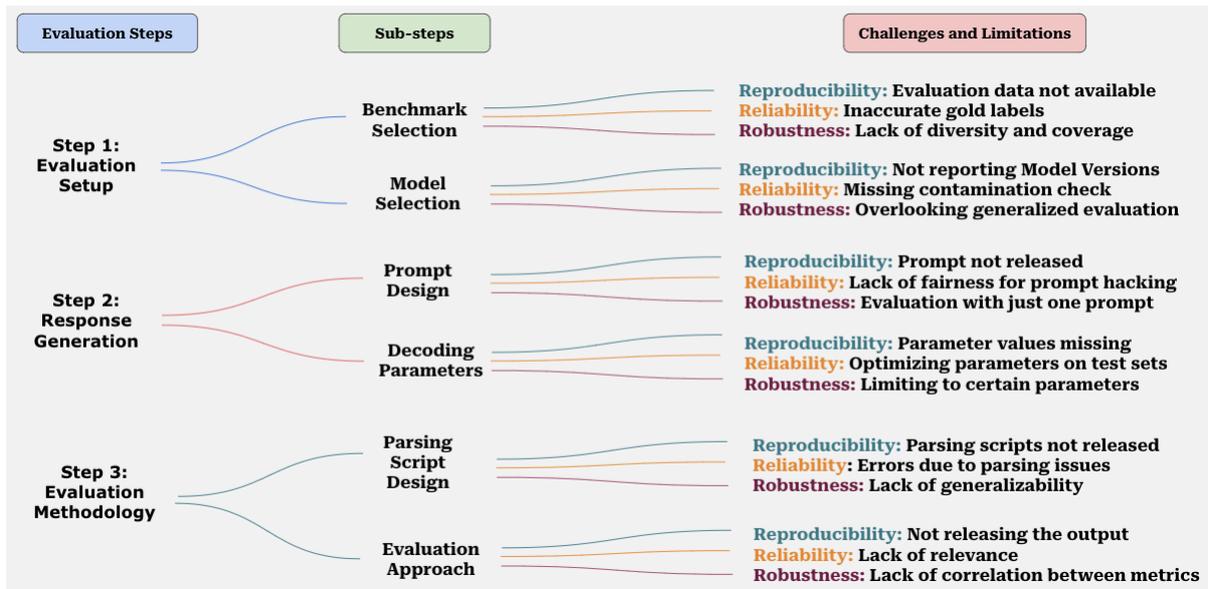


Figure 1: Typology of the LLM Evaluation Workflow. A more detailed description of the challenges and the limitations can be found in Table 5.

ducibility, reliability, and robustness (see Section 3). Based on our findings, we provide a principled guideline in Section 4 to address current limitations in LLM evaluation. The data and the code used in this paper are publicly available here: <https://github.com/ntunlp/Critical-Review-of-LLM-Eval>.

2 Overview of LLM Evaluation Process

The following components are crucial for LLM evaluation: *Evaluation Setup*, *Response Generation*, and *Evaluation Methodology* (Chang et al., 2024). Each component has its own challenges, which we discuss in Section 3. These components in an evaluation workflow are shown in Figure 1.

2.1 Evaluation Setup

Benchmark Selection: To initiate the evaluation process of LLMs, the first step is selecting appropriate benchmarks. We categorize the benchmarking datasets into the following: *general capability benchmarks*, *specialized benchmarks*, and *other diverse benchmarks*. We refer to general capability benchmarks as the ones that are often used for evaluation upon the release of an LLM (e.g., MMLU (Hendrycks et al., 2020b), HumanEval (Chen et al., 2021)). In addition, there are specialized benchmarks that measure specific capabilities of LLMs (e.g., MT-Bench for chatting capabilities (Zheng et al., 2024)). There are also other benchmarks that usually combine multiple benchmarks

to evaluate LLMs on diverse task (e.g., HELM (Liang et al., 2022)). We provide more details on each category in Appendix A.1.

Model Selection: Selecting the appropriate model from the numerous LLMs currently available is crucial for ensuring a fair evaluation, as it helps to avoid risks such as data contamination and unfair comparisons. For a detailed discussion on prominent LLMs, see Appendix A.2.

2.2 Response Generation

Once the benchmarks and the models are selected, the next step in the evaluation process is to design the prompt and set up the decoding parameters for response generation. In the *prompt design* step, decisions on what type of prompting (e.g., zero-shot or few-shot) would be used are taken. Moreover, configuring the *decoding parameters* (e.g., temperature) is important to ensure optimal performance (Shi et al., 2024). More discussions on this are provided in Appendix A.3 and A.4.

2.3 Evaluation Methodology

Parsing Script Design: Evaluating LLM-generated responses is difficult because they often produce verbose outputs (see Table 6 for some examples). Therefore, parsing scripts are often necessary (Jahan et al., 2024; Laskar et al., 2023a) to extract target labels before applying evaluation metrics, ensuring alignment with evaluation criteria to maintain reliability.

Availability (%)				Comparison (%)	
Prompt	Code	Prompt + Code	Model Version	Fair	Unfair
90.6	53.3	50.0	29.3	20.7	79.3

Table 1: Availability of resources and fairness in model comparisons (out of 212 papers), analyzed by Balloccu et al. (2024).

Evaluation Approach: The evaluation approach can be divided into the following: *automatic evaluation*, *human evaluation*, *LLMs as evaluators*. In *automatic evaluation*, before applying task-specific metrics (e.g., F1, Exact Match, Perplexity (Jelinek et al., 1977)), parsing scripts are often utilized to extract the targeted answer, especially in discriminative tasks. *Human evaluation* is required to ensure qualitative assessments of LLM responses (e.g., measuring clarity, coherence, factuality) (van der Lee et al., 2021). Recently, human evaluation based on the Elo-based rating system (Zheng et al., 2024) has gained a lot of attention. Since human evaluation is time-consuming, the utilization of *LLMs as evaluators* to assess other LLMs has become a popular evaluation approach (Chiang and Lee, 2023; Huang et al., 2024a). More details on LLM evaluation approaches are in Appendix A.6.1.

3 Challenges in Evaluating LLMs

We examine challenges and limitations in the evaluation process of LLMs based on three dimensions: *reproducibility*, *reliability*, and *robustness*.

3.1 Reproducibility

Reproducibility, the ability to consistently replicate model results under the same conditions, is a major challenge in generative models (Biderman et al., 2024). The primary challenge is the lack of comprehensive documentation for each part of the evaluation cycle, including benchmarking datasets, prompt construction, model details, decoding strategy, response parsing, and evaluation methodology (Kosch and Feger, 2024; McIntosh et al., 2024). Table 1 presents an analysis by Balloccu et al. (2024), revealing that a relatively low percentage of the analyzed papers shared their resources. Below, we discuss factors impacting reproducibility in the evaluation step.

3.1.1 Missing Details on Data & Models Used

Benchmarking Data: One factor that can negatively impede the ability to reproduce results is not

releasing the exact data used for evaluation (Balloccu et al., 2024). Many studies evaluate LLMs on only a subset of existing datasets (Bang et al., 2023; Kocoń et al., 2023), while others use the exact benchmarking datasets (Laskar et al., 2023a; Qin et al., 2023). Despite the expectation not to compare results across studies using different subsets of the data, such comparisons often occur, as discussed by Balloccu et al. (2024). Nonetheless, without explaining the sampling strategy, or releasing the subsets used for evaluation (and possibly their responses), reproducing results using different data subsets of the same size is challenging. **Model Versions:** The information regarding the version of a model being used is also missing in many studies (Balloccu et al., 2024; Biderman et al., 2024), creating reproducibility concern (see Table 1). The continuous updates of the closed-source models, often with undisclosed changes can also impact reproducibility. With these updates, earlier versions are often deprecated, and results from these versions may not apply to newer models (Chen et al., 2023b), making prior evaluation results to be no longer reproducible (Bang et al., 2023; Kocoń et al., 2023; Laskar et al., 2023a; Qin et al., 2023). Therefore, it is crucial to specify the model versions used (Balloccu et al., 2024; Biderman et al., 2024), while model owners should keep earlier versions available.

3.1.2 Lacking Response Generation Details

Prompting: The lack of details behind how the prompts are designed may make the findings in different literature inconsistent. For instance, variations in prompt design can lead to significantly different results, as seen in various studies (Bang et al., 2023; Jahan et al., 2024; Laskar et al., 2023a; Qin et al., 2023). While few-shot learning is found to outperform zero-shot in the original evaluation conducted by the authors of various LLMs (Anil et al., 2023; OpenAI, 2023; Touvron et al., 2023b), many independent evaluations demonstrate that adding few-shot examples does not necessarily outperform zero-shot models in every task (Jahan et al., 2024; Ye et al., 2023a). This raises the concern of whether certain prompt engineering techniques or optimizations to select few-shot samples were applied in the original evaluations. Hence, not disclosing the details behind how the prompt is designed or how the few-shot examples are selected can hinder reproducibility.

Decoding Strategy: LLMs are sensitive to de-

coding parameters, leading to significant performance variations based on the chosen settings (Roziere et al., 2023; Touvron et al., 2023b). However, crucial details on their selection are excluded in existing literature (Bang et al., 2023; Kocoń et al., 2023; Laskar et al., 2023a; OpenAI, 2023; Qin et al., 2023; Team et al., 2023). This lack of transparency raises reproducibility concerns, which could be responsible for inconsistent results across studies even when similar prompts are used. For instance, Qin et al. (2023) found that adding output length restrictions in the prompt to generate summaries in no more than N words led to a performance drop in the SAMSum dataset (Gliwa et al., 2019). However, Laskar et al. (2023a) found that such controlled experiments led to a gain in performance in the SAMSum dataset.

3.1.3 Evaluation Methods Unavailable

Parsing Scripts: LLM-generated responses often require parsing scripts to extract desired information. However, as demonstrated in Table 1, Balloccu et al. (2024) observed in their analysis that almost half of the LLM evaluation papers do not release any codes. We also observe that most studies (these include both the LLM technical reports, as well independent evaluations) do not release their parsing scripts (Bang et al., 2023; Kocoń et al., 2023; OpenAI, 2023; Qin et al., 2023; Team et al., 2023, 2024). Nonetheless, inaccurate design of parsing scripts may lead to different evaluation results (Laskar et al., 2023a). Thus, the unavailability of parsing scripts would complicate result comparisons while impacting reproducibility (Balloccu et al., 2024; Biderman et al., 2024).

Evaluation Approach: LLMs are increasingly used to evaluate other LLMs in development (Zheng et al., 2024). Concerns arise due to the use of closed-source LLMs as evaluators, as their frequent updates can affect reproducibility (Chen et al., 2023b; Verga et al., 2024). Moreover, Chen et al. (2023b) observed significant behavioral changes in closed-source LLMs over short periods. Such reproducibility concerns are also observed in prior research that used LLMs as evaluators. For instance, Chiang and Lee (2023); Zheng et al. (2024) found that using closed-source LLMs as the judge could collide with human evaluations, whereas Fu et al. (2023b) observed the opposite. Since the recently proposed Prometheus-2 (Kim et al., 2024a) model is an open-source alternative and demonstrates a strong correlation with

humans, utilizing open-source LLMs as the judge can help mitigate the reproducibility issues prevalent with closed-source LLMs.

3.2 Reliability

Reliability, the ability to trust that outcomes are as intended, is another challenge encountered during evaluation. Issues like contamination/inaccurate labels in the data, irrelevant evaluation methods, and unfair comparisons may impact the reliability of the findings, which we discuss below.

3.2.1 Data and Model Integrity Issues

Data Integrity: Errors in benchmarks undermine accurate conclusions and model comparisons, rendering evaluations of LLMs unreliable. An integrity-compromising factor is the presence of incorrect gold labels. For instance, existing issues in the gold labels of the widely used MMLU (Hendrycks et al., 2020b) dataset have led to the development of MMLU-Pro (Wang et al., 2024b) and MMLU-Redux (Gema et al., 2024). Recently it was also found that the coding benchmarks, HumanEval (Chen et al., 2021), lacked essential test cases, leading to the development of an advanced version, HumanEvalPlus (Liu et al., 2024b).

Despite these improvements, many recent studies continue to use the older versions of datasets. For instance, despite the release of HumanEvalPlus, HumanEval is still used to benchmark LLM coding performance (Gloeckle et al., 2024; Jiang et al., 2023; Li et al., 2023c; Roziere et al., 2023; Team et al., 2023, 2024; Wong et al., 2023), potentially providing misleading insights. In addition, outdated labels in existing benchmarks undermine reliability of gold references. For example, in tasks like open-domain question answering, which demand real-world knowledge, many gold labels become outdated over time, as noted by Laskar et al. (2023a). Consequently, even if LLMs produce correct answers, comparing them to obsolete gold labels can yield inaccurate results. Moreover, in tasks like summarization, LLM-generated summaries are often favored over human-annotated gold references (Ding et al., 2022; Pu et al., 2023; Zhang et al., 2024b).

Contamination in Existing Models: Contamination occurs when a benchmarking dataset is used in training, reducing result reliability and validity (Sainz et al., 2023; Shi et al., 2023; Zhou et al., 2023b). Ensuring benchmarking examples are excluded from training data is essential to maintain

reliable results. Since LLMs are pre-trained on vast amounts of text data available on the internet, this could lead to unfair evaluations if LLMs have already encountered these datasets during their pre-training phase (Balloccu et al., 2024; Ravaut et al., 2024; Xu et al., 2024).

Nonetheless, most prior LLM evaluation work focusing on zero-shot evaluation did not conduct any data contamination tests (Bang et al., 2023; Laskar et al., 2023a; OpenAI, 2023; Qin et al., 2023; Team et al., 2023), raising concerns about whether these evaluations truly represent the zero-shot capabilities of LLMs. Recent research has also demonstrated a strong possibility of data contamination in many datasets used to evaluate different LLMs (Balloccu et al., 2024; Golchin and Surdeanu, 2023; Li and Flanigan, 2023; Matton et al., 2024; Oren et al., 2023; Ravaut et al., 2024; Sainz et al., 2023; Xu et al., 2024; Zhang et al., 2024a). With the current generation of LLMs being extremely capable of learning new skills with minimal amounts of data, exposing them to evaluation data may undermine the measurement of their true capabilities. Since the possibility of data contamination has led to the development of new versions of existing datasets (e.g., utilizing GSM-8K to construct GSM-1K (Zhang et al., 2024a)), it is crucial to use fair evaluation datasets.

3.2.2 Lack of Fairness by Manipulating Response Generation

Prompt Hacking: One major concern in terms of lack of fairness in LLM evaluation is the possibility of prompt hacking (Schulhoff et al., 2023), which involves manipulating input prompts to a language model to elicit desired responses (e.g., biasing the outputs, or taking unfair advantages by using specific few-shot examples). While the performance of LLMs depends on many factors relevant to how the prompt is structured, most work (Bang et al., 2023; Laskar et al., 2023a; Qin et al., 2023), even the official technical reports (Anthropic, 2024; OpenAI, 2023; Team et al., 2023) of different LLMs lack the necessary details behind prompt construction (e.g., missing scientific validity on why a certain prompt was preferred over others, how the few-shot examples are selected, etc.). This makes the claims regarding the effectiveness and limitations of certain LLMs in comparison to others questionable². Recogniz-

²<https://crfm.stanford.edu/2024/05/01/helm-mm1u.html>

ing these parallels underscores the need for transparency and robust methodologies to ensure fairness in AI research and development.

Lack of Transparency in Decoding Parameters: Shi et al. (2024) demonstrated that extensive tuning of decoding parameters could improve the performance during inference. However, how the different decoding parameters are selected is often underexplored in existing evaluations (Bang et al., 2023; Laskar et al., 2023a,b; OpenAI, 2023; Qin et al., 2023; Team et al., 2023), as discussed in Section 3.1. This poses the risk of optimizing the parameters on test sets to improve performance.

3.2.3 Inappropriate Evaluation Methodology

Inaccurate Design of Parsing Scripts: As Laskar et al. (2023a) observed, evaluating LLMs entirely with an automated approach based on the answer extracted using parsing scripts may lead to an error of up to more than 10% difference in many tasks. This raises questions about the reliability of LLM evaluations that solely depend on parsing scripts without validating the scripts' effectiveness for the task. To tackle this, Laskar et al. (2023a) proposed a hybrid approach combining parsing script-based automatic evaluation with human-in-the-loop (Laskar et al., 2022a; Wu et al., 2022). Initially, the parsing script extracts answers from LLM-generated responses. If any issues arise, humans resolve them, enhancing the reliability of parsing-based automatic evaluation.

In Figure 2, we demonstrate the differences between automatic and hybrid evaluation in Open-Domain QA³ and reading comprehension datasets⁴. The figure highlights the influence of human intervention on results in open-domain QA, where LLMs may generate synonymous or time-sensitive correct answers, potentially rendering gold answers outdated (Laskar et al., 2023a). Parsing script-based automatic evaluation is found to be reliable in Race datasets for reading comprehension, whereas notable discrepancies are observed in the SQuAD-V2 dataset. Therefore, there's a need for designing dependable parsing scripts and involving humans when appropriate.

Evaluation Approaches Lacking Relevancy: In *generative tasks*, utilizing automatic string-based matching techniques may not be reliable as well.

³NQ-Open (Kwiatkowski et al., 2019), WebQuestions (Talmor and Berant, 2018), TriviaQA (Joshi et al., 2017))

⁴SQuAD-V2 (Rajpurkar et al., 2018), Race-High and Race-Middle (Lai et al., 2017)

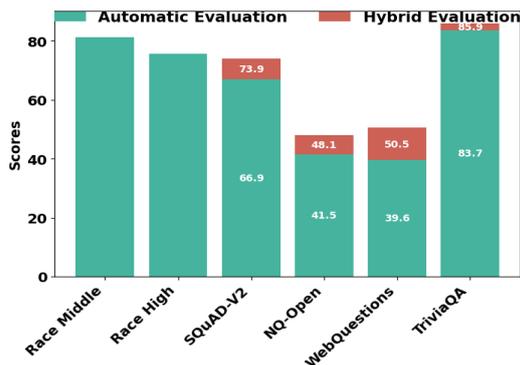


Figure 2: Comparing Automatic and Hybrid Evaluation.

For instance, [Laskar et al. \(2023a\)](#) observed that despite LLMs scoring quite poorly on the ROUGE metric compared to SOTA summarization models, humans often prefer LLM-generated responses. Moreover, recent research observed potential biases while using LLMs as evaluators, such as LLMs preferring responses generated by LLMs of the same series, positional bias ([Bai et al., 2024](#); [Stureborg et al., 2024](#); [Wang et al., 2023b](#); [Wu and Aji, 2023](#)). To mitigate this, [Verga et al. \(2024\)](#) proposed a new technique that leveraged multiple LLMs as juries instead of using a single LLM as the judge. This approach demonstrates higher correlations with humans, while mitigating biases.

3.3 Robustness

In the context of evaluating LLMs, robustness refers to the model’s ability to maintain consistent performance across a wide range of inputs, conditions, or tasks. While there are many evaluation benchmarks currently available, existing work mostly relies on evaluating LLMs on some common benchmarks. This raises the question of whether the performance of LLMs in these common benchmarks reflects their true capabilities and limitations. In this section, we study the robustness of existing LLM evaluations.

3.3.1 Lacking Generalized Evaluation

Limiting Evaluation to Certain Scenarios: Interestingly, it has been observed in recent research that certain performance gains in a specific dataset may not necessarily imply that it would also improve the performance in other datasets for similar tasks ([Jahan et al., 2024](#); [SambaNova, 2024](#)). For instance, [Jahan et al. \(2024\)](#) observes that not a single LLM has superiority over other LLMs

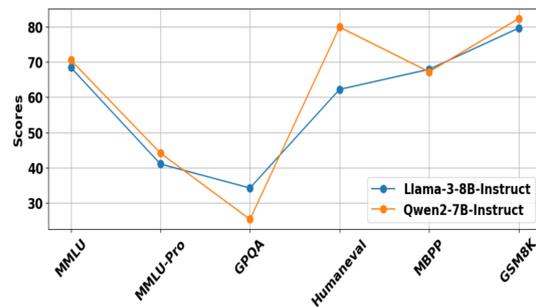


Figure 3: Performance Comparison: LLaMA-3 and Qwen2

Tokenizer	Vocab	MMLU	MMLU-Pro	MixEval	MixEval-Hard
LLaMA-2	32,000	0.52	0.45	0.29	0.11
LLaMA-3	128,256	0.27	0.21	0.09	0.03
Mistral	32,000	0.59	0.51	0.31	0.11
Qwen2	151,646	0.22	0.17	0.08	0.02

Table 2: Comparison of vocabulary coverage across different datasets and LLM tokenizers. The scores represent the percentage of tokenizer vocabulary that is covered by the respective dataset.

across all biomedical datasets and tasks. This is also evident if we compare the results between LLaMA-3 and Qwen2 reported in ([Qwen2, 2024](#); [Yang et al., 2024](#)). As shown in Figure 3, while the Qwen2 model outperforms LLaMA-3 on most datasets, it falls short on GPQA and MBPP. Interestingly, for coding tasks, Qwen2 significantly outperforms LLaMA-3 on the HumanEval dataset ([Chen et al., 2021](#)) but not on the MBPP dataset ([Austin et al., 2021](#)). Meanwhile, existing common benchmarks also do not take into account some specific settings, such as how LLMs perform in long context scenarios, as recent research demonstrated that LLMs often struggle to generate the correct answer when relevant information does not appear at the beginning or end of the input context ([Liu et al., 2024c](#)). This highlights the importance of evaluating the generalized performance of LLMs across a set of diverse benchmarks and settings, instead of limiting evaluation to only common benchmarks like MMLU ([Hendrycks et al., 2020b](#)).

Diversity and Coverage in Benchmarks: Although benchmarking datasets are designed to address specific problems and objectives, the variation and complexity of language within these datasets are often unclear. [Liang et al. \(2022\)](#) highlighted that better coverage in benchmarking datasets would enhance the comprehensiveness of the model’s evaluation. While different language

models use different tokenizers to represent the benchmarking dataset, it also leads to variations in what is evaluated across models.

As can be seen in Table 2, we conducted a small-scale analysis for LLaMA-2 (Touvron et al., 2023b), LLaMA-3,⁵ Mistral (Jiang et al., 2023), and Qwen2⁶ on two benchmarking datasets with varying complexities: MMLU (Hendrycks et al., 2020b) and its more challenging version, MMLU-Pro (Wang et al., 2024b), as well as MixEval (Ni et al., 2024) and its harder version, MixEval-Hard. Our findings indicate that these datasets cover a relatively small portion of the model’s capabilities. Specifically, for MixEval, as the datasets became more diverse and dynamic, the vocabulary coverage for the tokenizer decreased. This trend continued as the datasets increased in difficulty, with vocabulary coverage further declining.

3.3.2 No Tuning of Prompt and Decoding Parameters

While various combinations of decoding parameters may lead to differences in results (Shi et al., 2024), possibly due to high computing requirements, existing LLM evaluation work mostly undermines the necessity of evaluating how the model performance may vary depending on its variations. Similar to the absence of decoder parameter tuning, most prior work also evaluated LLMs using only a single prompt (Bang et al., 2023; Jahan et al., 2024; Kocoń et al., 2023; Laskar et al., 2023a; Qin et al., 2023). However, in the real world, users express themselves with diverse word choices, varying semantics and syntaxes, alongside minor discrepancies (e.g., misspellings or differing punctuation styles). To further examine the effects of prompt variations, we conduct an experiment using GPT-4o (2024-04-09) and GPT-3.5-Turbo (0125) (OpenAI, 2023), as well as Claude-3-Opus (2024-02-29) (Anthropic, 2024) with the prompts used by (Laskar et al., 2023a) and (Qin et al., 2023) in the SAMSum dataset. For this experiment, the default parameters for respective LLMs are used.

As shown in Figure 4, the restricted prompting method by Laskar et al. (2023a) consistently outperforms the unrestricted approach across all three models. Conversely, the restricted prompting method by Qin et al. (2023) fails to surpass

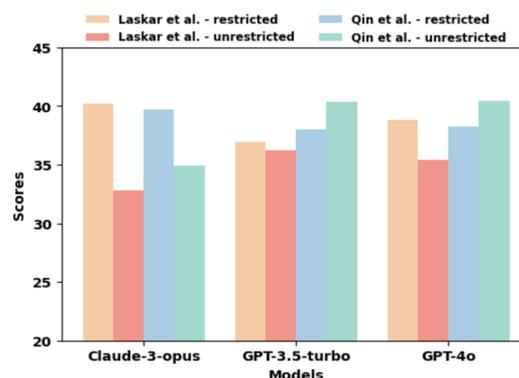


Figure 4: ROUGE-1 scores in the SAMSum dataset based on Prompt Tuning.

the unrestricted approach for GPT-3.5 and GPT-4o. However, it surprisingly outperforms the unrestricted method, indicating the significant impact of prompt tuning across models. Evaluating language models with a single prompt lacks fairness (Zhu et al., 2023b), yet it remains common practice (Bang et al., 2023; Laskar et al., 2023a; Qin et al., 2023). Minor prompt variations can lead to diverse outcomes for different models (Alzahrani et al., 2024; An et al., 2023; Biderman et al., 2024; Lanham et al., 2023; Sclar et al., 2023; Wei et al., 2024; Zhang et al., 2024a), highlighting the need to compare benchmarks across multiple prompts. Using automated prompt tuning techniques like Meta Probing Agents (Zhu et al., 2024) can ensure robustness to prompt variations.

3.3.3 Evaluation Method’s Generalizability and Correlation Shortcomings

While automatic evaluations are usually utilized in discriminative tasks, they may not be applicable to every task, as demonstrated by Jahan et al. (2024) that parsing scripts are not usable in certain discriminative tasks like relation extraction. Jahan et al. (2024) also noted a significant performance gap between the string-matching-based ROUGE metric (Lin, 2004) and the contextual similarity-based metric BERTScore (Zhang et al., 2019) in text summarization. While larger models achieve better accuracy, they involve a speed-accuracy trade-off (Parvez et al., 2019), leading to higher costs and latency (Fu et al., 2024b; Laskar et al., 2023b). While metrics like perplexity are widely used to evaluate language models (Chen et al., 2023c), Huang et al. (2024b) found that quantized LLaMA-3 versions have lower output confidence than the original. They noted simi-

⁵<https://llama.meta.com/llama3/>

⁶<https://github.com/QwenLM/Qwen2>

Model	Chatbot Arena	HELM MMLU	Vellum MMLU
GPT-4o-2024-05-13	1 (1)	2 (2)	1 (1)
GPT-4-Turbo-2024-04-09	5 (3)	3 (3)	3 (3)
GPT-4-0125-preview	6 (4)	5 (5)	4 (4)
Gemini-1.5-Pro	4 (2)	4 (4)	13 (6)
Gemini-1.5-Flash	10 (6)	10 (6)	10 (5)
Claude-3-Opus-2024-02-29	7 (5)	1 (1)	2 (2)

Table 3: Rankings of models on LMSys Chatbot Arena vs two MMLU implementations. The relative rank of each model in MMLU is shown in parentheses.

lar model rankings for perplexity and a common-sense QA dataset. However, [Hu et al. \(2024\)](#) found no correlation between perplexity and long context understanding tasks, highlighting the need for robust evaluations with human-correlated metrics.

This raises another question, whether automated evaluations and LLM-as-a-judge correlate with human evaluations (e.g., Elo ratings). [Zheng et al. \(2024\)](#) demonstrated significant correlations between Elo ratings, LLM-as-a-judge, and automated evaluations. However, recent research ([Alzahrani et al., 2024](#)) suggest that automated evaluations, especially those using multiple-choice questions, can yield unstable rankings with minor changes in evaluation methods. Given this instability, it prompts us to question why these automated tests should align with human Elo ratings despite demonstrating such inconsistencies. In our view, we should focus not only on correlating scores but also on how well a benchmark’s rankings align with the gold standards. Analysis in Table 4 for GPT-4 ([OpenAI, 2023](#)), Gemini ([Team et al., 2023](#)), and Claude-3 ([Anthropic, 2024](#)) reveals two key observations: (i) MMLU rankings disagree with LMSys Chatbot Arena and (ii) MMLU rankings vary among themselves due to implementation differences.

4 Recommendations and Best Practices

So far, we’ve outlined the primary challenges in evaluating LLMs. In light of these challenges, a crucial question arises: How can we enhance the evaluation of LLMs? Crafting a structured framework that’s both practical and easy to implement is daunting, given the complexities of generative LLM development. Previous studies tended to focus on specific evaluation aspects without offering comprehensive guidelines for the entire evaluation cycle, leaving researchers without clear guidance. Before diving into recommendations for each evaluation stage, it’s important to acknowledge three

key factors shaping current LLM evaluation practices: inherent randomness in generative models, significant computational demands, and insufficient documentation across stages.

Evaluation Setup: Selecting benchmarks for model assessment is crucial. Rather than simply replicating past choices, researchers should align datasets with required capabilities. To ensure *robustness*, datasets should vary across expected LLM capabilities (e.g., long-context understanding), tasks (e.g., summarization), and language complexity (e.g., vocabulary coverage). Ideally, a metric should measure dataset diversity. For model selection, conduct contamination tests between the chosen model and benchmarks using relevant techniques ([Ravaut et al., 2024](#)). This acts as an additional filter for benchmarking datasets, ensuring selection of unseen ones measuring intended capabilities. Meanwhile, for *reproducibility*, document any subset use of benchmarking datasets, along with the selected model version. In addition, throughout scientific history, intelligence progress has evolved across generations. Tests from a decade ago may appear simplistic compared to today’s standards (e.g., Math Olympiads, ICPC programming contests). Refreshing LLM evaluations periodically can effectively communicate standard capabilities in both open and closed-source LLM markets and ecosystems (e.g., chatbots). Hence, to ensure *reliability*, verify if the dataset has updated versions and incorporate them if available (e.g., HumanEvalPlus ([Liu et al., 2024b](#)), MMLU-Pro ([Wang et al., 2024b](#)), GSM-1K ([Zhang et al., 2024a](#)))

Response Generation: For *reproducibility*, thorough documentation of prompts (e.g., explaining the selection of few-shot samples) and parameter settings (e.g., use tools like mlflow⁷ or Weights & Biases⁸ (W&B)) is essential. To ensure *reliability*, it’s crucial to justify why specific prompts and parameters are chosen over others by providing comparisons with alternative options. As for *robustness*, experimenting with diverse prompts and parameters is the key to showcasing their effectiveness and limitations in different scenarios. In resource-constrained environments, conducting experiments with diverse evaluation settings may pose challenges, yet it remains vital to perform robust evaluations on at least a subset of samples.

⁷<https://mlflow.org/>

⁸<https://wandb.ai/site>

Step	Sub-Step	Recommendation	Implementation: Suggested Tools or Techniques
Evaluation Setup	Benchmark Selection	Selected benchmarks should align with the capabilities required and updated versions of the datasets should be used to ensure reliability, diversity in the selected benchmarks is required to ensure robustness, and proper documentation of the dataset subsets is required for reproducibility.	Reliability: Use refined benchmarks like MMLU-Pro, Human-Eval Plus, GSM-1k to address the limitations in existing benchmarks to improve reliability. Reproducibility: Document the data sampling technique and release the data subset used for evaluation alongside the model-generated response. Robustness: Check tokenizer vocabulary coverage in selected benchmarks.
	Model Selection	Data contamination check in the selected model is required for reliability, proper versioning of the model is required for reproducibility, and diverse capability evaluation (e.g., latency, memory usage, format following capability, etc.) is important to ensure robustness.	Reliability: Use tools like LLMsSanitize Library (Ravaut et al., 2024) for contamination check. Reproducibility: Use mflow or W&B for documentation. Robustness: Use tools like pyNVML to measure GPU memory requirements, FOFO for format following, compare accuracy vs latency trade-off, etc.
Response Generation	Prompt Design	Release the prompts and few-shot examples for reproducibility, justify the selection of certain prompts and few-shot examples to ensure reliability, and compare with alternative prompts to ensure robustness.	Reliability: Justify the choice of certain prompts to ensure no potential of prompt hacking and compare the alternatives. Meanwhile, clearly demonstrate what and how few-shot examples are selected. Reproducibility: Use tools like LM-Evaluation-Harness. Robustness: Use Prompt Bench or Meta-Probing Agent.
	Decoding Parameters	Document the decoding parameters to ensure reproducibility, justify the selection to ensure reliability, and experiment with various parameters to ensure robustness.	Reliability: Justify the choice of certain parameters to eliminate the risk of optimization in the test data. Reproducibility: Use mflow or W&B. Robustness: Compare the performance based on different decoding parameters, at least in a subset of the data.
Evaluation Methodology	Parsing Script Design	Accurate parsing of the response is required for reliability, availability of these scripts is needed for reproducibility, and parsing scripts should show robustness across different models and datasets.	Reliability: Validate the reliability based on human evaluation, at least on a subset. Reproducibility: Release the code. Robustness: Evaluate multiple models and datasets, across all types of labels and corner cases.
	Evaluation Approach	Availability of the evaluation output is required for reproducibility, selected evaluation metrics should maintain correlation with humans to ensure reliability, and multiple evaluation metrics are required for evaluation robustness.	Reliability: Validate the effectiveness of selected metrics (e.g., measure correlation with humans), use techniques like LLM-as-juries to mitigate bias. Reproducibility: Release the Evaluation Output. Robustness: Use multiple evaluation metrics (e.g., in Summarization, use both word-based (e.g., ROUGE) or Contextualized (e.g., BERTScore) metrics), measure latency, GPU usage via pyNVML.

Table 4: Recommendations and Implementation Suggestions.

Evaluation Methodology: To ensure *reproducibility*, the parsing scripts and the output data used for evaluation should be published. Meanwhile, sanity-checking on the parsing script should be done to ensure *reliability* and *robustness* of the designed parsing script. This can be done by creating test cases for various response types, and then verifying (with human intervention if possible) whether the parsing script can reliably extract the targeted answer from the generated response. Meanwhile, reliance on string-based metrics like ROUGE should be minimized in favor of qualitative evaluations to ensure the *reliability* of the chosen evaluation methodology. Given the cost and time constraints of human qualitative evaluation, LLM-based evaluators can be used as alternatives but must be validated for potential biases (e.g., multiple LLMs as juries instead of using a single LLM as the judge (Zheng et al., 2024)). Finally, *robust* evaluation using task-specific metrics is encouraged. For this purpose, metrics that lack alignment with humans should be avoided. Moreover, measuring runtime latency using tools like

pyNVML⁹ is recommended to evaluate the real-world applicability of different LLMs.

5 Conclusions and Future Work

In this paper, we systematically survey the challenges and limitations in evaluating LLMs. We identified significant inconsistencies and complexities at various stages of the evaluation pipeline, impacting the reproducibility, reliability, and robustness of the results. These issues underline the necessity for a standardized and systematic approach for LLM evaluation to ensure their reliable usage in real-world applications. By comprehensively reviewing the current evaluation practices, we have provided a set of recommendations aimed at enhancing the consistency and fairness of LLM evaluations. Therefore, future work should focus on developing and adopting standardized evaluation protocols for LLMs to address the identified limitations. This includes creating benchmark datasets, evaluation metrics, and proper documentation of the evaluation settings to ensure reproducibility, reliability, and robustness.

⁹<https://pypi.org/project/pynvml/>

Acknowledgements

We would like to thank all the anonymous reviewers for their excellent review comments. This research was supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada and the York Research Chairs (YRC) program. We also acknowledge Compute Canada for the computing resources. Finally, we thank Mir Tafseer Nayeem for providing valuable feedback.

Limitations

One limitation of this work is that it is focused only on the evaluation phase of the LLM development cycle. Therefore, the challenges and limitations that happen during the training phase of LLMs are left out of the scope of this paper. Nonetheless, with the rapid growth of LLM technologies and huge financial incentives, it is essential to conduct a fair and reliable evaluation of LLM, alongside ensuring robustness and reproducibility, which is the focus of this work.

Another limitation of this study is that it does not study how to prevent closed-source LLMs from getting access to the online benchmarks. For instance, assume we have two entities: model developers and evaluators. Evaluators do not want to expose their data to the modeling team. Conversely, model developers do not want to release their model weights due to significant financial incentives. If evaluators use an API to get the responses, there is a risk that the queries may get exposed to the model developers. Therefore, without getting access to the weights, evaluators cannot reliably assess the models on their queries. Mathematically and technically, there is no fundamental way to solve this problem without altering the training dynamics which may not be an option for training teams.

Moreover, given the limited amount of study to evaluate LLMs in non-English data, our work was more focused on the monolingual scenario (mostly on English data). Therefore, investigating the challenges and limitations of LLM evaluation in multilingual and resource-constrained scenarios could be studied in the future, alongside also studying the performance of various tokenizers (both multilingual and monolingual) in LLM benchmarking (Choo and Kim, 2023; Rust et al., 2021)).

Finally, the multimodal capability, in other words, the ability to understand both language and

vision is another interesting capability of recently proposed LLMs (Bai et al., 2023; Chen et al., 2023a; Dai et al., 2024; Liu et al., 2023b, 2024a; Luo et al., 2024; Ye et al., 2023b; Zhang et al., 2023; Zhu et al., 2023a). This has led to the development of many multi-modal benchmarks (Chen et al., 2024b; Fu et al., 2023a, 2024a; Guan et al., 2023; Li et al., 2023a,b,d; Liu et al., 2024a, 2023e; Lu et al., 2022; Qiu et al., 2024; Yu et al., 2023). However, this paper was mostly focused on text-based NLP tasks and the evaluation of LLMs on multimodal benchmarks is left out for future work.

Ethics Statement

This paper only reviews the existing challenges and limitations in LLM evaluations and provides an opinion piece and recommendation to ensure reliable, robust, and reproducible evaluations of LLMs. Thus, this review does not pose any ethical concerns.

References

- Ahmed Abdelali, Hamdy Mubarak, Shammur Chowdhury, Maram Hasanain, Basel Mousi, Sabri Boughorbel, Samir Abdaljalil, Yassine El Kheir, Daniel Izham, Fahim Dalvi, et al. 2024. Larabench: Benchmarking arabic ai with large language models. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 487–520.
- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, et al. 2023. Mega: Multilingual evaluation of generative ai. *arXiv preprint arXiv:2303.12528*.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay,

- Quentin Malartic, et al. 2023. The falcon series of open language models. *arXiv preprint arXiv:2311.16867*.
- Iñigo Alonso, Maite Oronoz, and Rodrigo Agerri. 2024. Medexpqa: Multilingual benchmarking of large language models for medical question answering. *arXiv preprint arXiv:2404.05590*.
- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M Saiful Bari, and Haidar Khan. 2024. [When benchmarks are targets: Revealing the sensitivity of large language model leaderboards](#).
- Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Weizhu Chen, and Jian-Guang Lou. 2023. Skill-based few-shot selection for in-context learning. *arXiv preprint arXiv:2305.14210*.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Anthropic. 2024. [The claude 3 model family: Opus, sonnet, haiku](#).
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.
- Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, et al. 2024. Benchmarking foundation models with language-model-as-an-examiner. *Advances in Neural Information Processing Systems*, 36.
- Simone Balloccu, Patrícia Schmidová, Mateusz Lango, and Ondřej Dušek. 2024. Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source llms. *arXiv preprint arXiv:2402.03927*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
- Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Alham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive, Anthony DiPofi, Julen Etxaniz, Benjamin Fattori, Jessica Zosa Forde, Charles Foster, Mimansa Jaiswal, Wilson Y. Lee, Haonan Li, Charles Lovering, Niklas Muennighoff, Ellie Pavlick, Jason Phang, Aviya Skowron, Samson Tan, Xiangru Tang, Kevin A. Wang, Genta Indra Winata, François Yvon, and Andy Zou. 2024. [Lessons from the trenches on reproducible evaluation of language models](#).
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Meriem Boubdir, Edward Kim, Beyza Ermis, Sara Hooker, and Marzieh Fadaee. 2023. Elo uncovered: Robustness and best practices in language model evaluation. *arXiv preprint arXiv:2311.17295*.
- Sabri Boughorbel, MD Parvez, and Majd Hawasly. 2024. Improving language models trained with translated data via continual pre-training and dictionary learning analysis. *arXiv preprint arXiv:2405.14277*.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Jannis Bulian, Christian Buck, Wojciech Gajewski, Benjamin Boerschinger, and Tal Schuster.

2022. Tomayto, tomahto. beyond token-level answer equivalence for question answering evaluation. *arXiv preprint arXiv:2202.07654*.
- Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. 2024. A survey on mixture of experts. *arXiv preprint arXiv:2407.06204*.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. [A survey on evaluation of large language models](#). *ACM Trans. Intell. Syst. Technol.*, 15(3).
- Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwal, Edgar Dobriban, Nicolas Flammarion, George J Pappas, Florian Tramèr, et al. 2024. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. *arXiv preprint arXiv:2404.01318*.
- Anthony Chen, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2020. Mocha: A dataset for training and evaluating generative reading comprehension metrics. *arXiv preprint arXiv:2010.03636*.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024a. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. 2024b. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023a. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*.
- Lingjiao Chen, Matei Zaharia, and James Zou. 2023b. How is ChatGPT’s behavior changing over time? *arXiv preprint arXiv:2307.09009*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023c. Longlora: Efficient fine-tuning of long-context large language models. *arXiv preprint arXiv:2309.12307*.
- Steffi Chern, Ethan Chern, Graham Neubig, and Pengfei Liu. 2024. Can large language models be trusted for evaluation? scalable meta-evaluation of llms as evaluators via agent debate. *arXiv preprint arXiv:2401.16788*.
- Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937*.
- Sanghyun Choo and Wonjoon Kim. 2023. A study on the evaluation of tokenizer performance in natural language processing. *Applied Artificial Intelligence*, 37(1):2175112.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *arXiv preprint arXiv:2210.11416*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages. *arXiv preprint arXiv:2003.05002*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.

- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2024. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36.
- Fahim Dalvi, Maram Hasanain, Sabri Boughorbel, Basel Mousi, Samir Abdaljalil, Nizi Nazar, Ahmed Abdelali, Shammur Absar Chowdhury, Hamdy Mubarak, Ahmed Ali, et al. 2023. Llmehub: A flexible framework for accelerating llms benchmarking. *arXiv preprint arXiv:2308.04945*.
- Bosheng Ding, Chengwei Qin, Linlin Liu, Lidong Bing, Shafiq Joty, and Boyang Li. 2022. Is gpt-3 a good data annotator? *arXiv preprint arXiv:2212.10450*.
- Markus Freitag and Yaser Al-Onaizan. 2017. Beam search strategies for neural machine translation. *arXiv preprint arXiv:1702.01806*.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiwu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2023a. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.
- Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. 2024a. Blink: Multimodal large language models can see but not perceive. *arXiv preprint arXiv:2404.12390*.
- Xue-Yong Fu, Md Tahmid Rahman Laskar, Cheng Chen, and Shashi Bhushan Tn. 2023b. [Are large language models reliable judges? a study on the factuality evaluation capabilities of LLMs](#). In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 310–316, Singapore. Association for Computational Linguistics.
- Xue-Yong Fu, Md Tahmid Rahman Laskar, Elena Khasanova, Cheng Chen, and Shashi Bhushan TN. 2024b. Tiny titans: Can smaller large language models punch above their weight in the real world for meeting summarization? *arXiv preprint arXiv:2402.00841*.
- Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. 2023a. Human-like summarization evaluation with ChatGPT. *arXiv preprint arXiv:2304.02554*.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023b. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, et al. 2024. Are we done with mmlu? *arXiv preprint arXiv:2406.04127*.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsun corpus: A human-annotated dialogue dataset for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79.
- Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. 2024. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*.
- Shahriar Golchin and Mihai Surdeanu. 2023. Time travel in llms: Tracing data contamination in large language models. *arXiv preprint arXiv:2308.08493*.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. 2023. Hallusionbench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566*.
- Yue Guo, Zian Xu, and Yi Yang. 2023a. Is ChatGPT a financial expert? evaluating language models on financial natural language processing. *arXiv preprint arXiv:2310.12664*.
- Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, Deyi Xiong, et al. 2023b. Evaluating large language models:

- A comprehensive survey. *arXiv preprint arXiv:2310.19736*.
- Rishav Hada, Varun Gumma, Adrian de Wynter, Harshita Diddee, Mohamed Ahmed, Monojit Choudhury, Kalika Bali, and Sunayana Sitaram. 2023. Are large language model-based evaluators the solution to scaling up multilingual evaluation? *arXiv preprint arXiv:2309.07462*.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. 2021. [Measuring coding challenge competence with APPS](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2020a. Aligning ai with shared human values. *arXiv preprint arXiv:2008.02275*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020b. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Yutong Hu, Quzhe Huang, Mingxu Tao, Chen Zhang, and Yansong Feng. 2024. Can perplexity reflect large language model’s ability in long text understanding? *arXiv preprint arXiv:2405.06105*.
- Hui Huang, Yingqi Qu, Jing Liu, Muyun Yang, and Tiejun Zhao. 2024a. An empirical study of llm-as-a-judge for llm evaluation: Fine-tuned judge models are task-specific classifiers. *arXiv preprint arXiv:2403.02839*.
- Wei Huang, Xudong Ma, Haotong Qin, Xingyu Zheng, Chengtao Lv, Hong Chen, Jie Luo, Xiaojuan Qi, Xianglong Liu, and Michele Magno. 2024b. How good are low-bit quantized llama3 models? an empirical study. *arXiv preprint arXiv:2404.14047*.
- Md Ashraf Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. 2024a. Mapcoder: Multi-agent code generation for competitive problem solving. *arXiv preprint arXiv:2405.11403*.
- Mohammed Saidul Islam, Raian Rahman, Ahmed Masry, Md Tahmid Rahman Laskar, Mir Tafseer Nayeem, and Enamul Hoque. 2024b. Are large vision language models up to the challenge of chart comprehension and reasoning? an extensive investigation into the capabilities and limitations of llms. *arXiv preprint arXiv:2406.00257*.
- Israt Jahan, Md Tahmid Rahman Laskar, Chun Peng, and Jimmy Huang. 2023. Evaluation of ChatGPT on biomedical tasks: A zero-shot comparison with fine-tuned generative transformers. In *The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, pages 326–336, Toronto, Canada. Association for Computational Linguistics.
- Israt Jahan, Md Tahmid Rahman Laskar, Chun Peng, and Jimmy Xiangji Huang. 2024. A comprehensive evaluation of large language models on benchmark biomedical text processing tasks. *Computers in Biology and Medicine*, page 108189.
- Fred Jelinek, Robert L Mercer, Lalit R Bahl, and James K Baker. 1977. Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America*, 62(S1):S63–S63.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611.
- Mohsinul Kabir, Mohammed Saidul Islam, Md Tahmid Rahman Laskar, Mir Tafseer Nayeem, M Saiful Bari, and Enamul Hoque. 2023.

- Benllmeval: A comprehensive evaluation into the potentials and pitfalls of large language models on bengali NLP. *arXiv preprint arXiv:2309.13173*.
- Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 235–251. Springer.
- Zachary Kenton, Noah Y Siegel, János Kramár, Jonah Brown-Cohen, Samuel Albanie, Jannis Bulian, Rishabh Agarwal, David Lindner, Yunhao Tang, Noah D Goodman, et al. 2024. On scalable oversight with weak llms judging strong llms. *arXiv preprint arXiv:2407.04622*.
- Mohammad Abdullah Matin Khan, M Saiful Bari, Xuan Long Do, Weishi Wang, Md Rizwan Parvez, and Shafiq Joty. 2023. xcodeeval: A large scale multilingual multitask benchmark for code understanding, generation, translation and retrieval. *arXiv preprint arXiv:2303.03004*.
- Md Tawkat Islam Khondaker, Abdul Waheed, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. Gptaraeval: a comprehensive evaluation of ChatGPT on arabic NLP. *arXiv preprint arXiv:2305.14976*.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024a. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*.
- Tae Soo Kim, Yoonjoo Lee, Jamin Shin, Young-Ho Kim, and Juho Kim. 2024b. Evallm: Interactive evaluation of large language model prompts on user-defined criteria. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–21.
- Masamune Kobayashi, Masato Mita, and Mamoru Komachi. 2024. Large language models are state-of-the-art evaluator for grammatical error correction. *arXiv preprint arXiv:2403.17540*.
- Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. *arXiv preprint arXiv:2302.14520*.
- Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielanievicz, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, et al. 2023. Chatgpt: Jack of all trades, master of none. *Information Fusion*, 99:101861.
- Thomas Kosch and Sebastian Feger. 2024. Risk or chance? large language models and reproducibility in human-computer interaction research. *arXiv preprint arXiv:2404.15782*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: Large-scale reading comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*.
- Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023. ChatGPT beyond english: Towards a comprehensive evaluation of large language models in multilingual learning. *arXiv preprint arXiv:2304.05613*.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*.
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. 2023. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*.
- Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Huang. 2023a. A systematic study and comprehensive evaluation of

- ChatGPT on benchmark datasets. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 431–469, Toronto, Canada. Association for Computational Linguistics.
- Md Tahmid Rahman Laskar, Cheng Chen, Xueyong Fu, and Shashi Bhushan Tn. 2022a. Improving named entity recognition in telephone conversations via effective active learning with human in the loop. *arXiv preprint arXiv:2211.01354*.
- Md Tahmid Rahman Laskar, Xue-Yong Fu, Cheng Chen, and Shashi Bhushan Tn. 2023b. Building real-world meeting summarization systems using large language models: A practical perspective. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 343–352.
- Md Tahmid Rahman Laskar, Enamul Hoque, and Jimmy Xiangji Huang. 2022b. Domain adaptation with pre-trained transformers for query-focused abstractive text summarization. *Computational Linguistics*, 48(2):279–320.
- Md Tahmid Rahman Laskar, Xiangji Huang, and Enamul Hoque. 2020. Contextualized embeddings based transformer encoder for sentence similarity modeling in answer selection task. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5505–5514.
- Md Tahmid Rahman Laskar, Mizanur Rahman, Israt Jahan, Enamul Hoque, and Jimmy Huang. 2023c. Can large language models fix data annotation errors? an empirical study using debatepedia for query-focused text summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10245–10255.
- Md Tahmid Rahman Laskar, Mizanur Rahman, Israt Jahan, Enamul Hoque, and Jimmy Huang. 2023d. CQSumDP: a ChatGPT-annotated resource for query-focused abstractive summarization based on debatepedia. *arXiv preprint arXiv:2305.06147*.
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Kraemer. 2021. Human evaluation of automatically generated text: Current trends and best practice guidelines. *Computer Speech & Language*, 67:101151.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2023a. Seed-bench-2: Benchmarking multimodal large language models. *arXiv preprint arXiv:2311.17092*.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023b. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*.
- Changmao Li and Jeffrey Flanigan. 2023. Task contamination: Language models may not be few-shot anymore. *arXiv preprint arXiv:2312.16337*.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023c. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023d. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. 2023e. Large language models in finance: A survey. In *Proceedings of the Fourth ACM International Conference on AI in Finance*, pages 374–382.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland

- Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022. [Competition-level code generation with alphacode](#). *Science*, 378(6624):1092–1097.
- Zongxia Li, Ishani Mondal, Yijun Liang, Huy Nghiem, and Jordan Lee Boyd-Graber. 2024. [Pedants \(precise evaluations of diverse answer nominee text for skinflints\): Efficient evaluation analysis and benchmarking for open-domain question answering](#).
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soyulu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. [Holistic evaluation of language models](#). *arXiv preprint arXiv:2211.09110*.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. [Truthfulqa: Measuring how models mimic human falsehoods](#). *arXiv preprint arXiv:2109.07958*.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023a. [Mitigating hallucination in large multi-modal models via robust instruction tuning](#). In *The Twelfth International Conference on Learning Representations*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023b. [Improved baselines with visual instruction tuning](#). *arXiv preprint arXiv:2310.03744*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024a. [Visual instruction tuning](#). *Advances in neural information processing systems*, 36.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2024b. [Is your code generated by ChatGPT really correct? rigorous evaluation of large language models for code generation](#). *Advances in Neural Information Processing Systems*, 36.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024c. [Lost in the middle: How language models use long contexts](#). *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Zhuoer Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, et al. 2023c. [Alignbench: Benchmarking chinese alignment of large language models](#). *arXiv preprint arXiv:2311.18743*.
- Yi Liu, Lianzhe Huang, Shicheng Li, Sishuo Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023d. [Recall: A benchmark for llms robustness against external counterfactual knowledge](#). *arXiv preprint arXiv:2311.08147*.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023e. [Mmbench: Is your multi-modal model an all-around player?](#) *arXiv preprint arXiv:2307.06281*.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. [Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts](#). *arXiv preprint arXiv:2310.02255*.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. [Learn to explain: Multimodal reasoning via thought chains for science question answering](#). *Advances in Neural Information Processing Systems*, 35:2507–2521.
- Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. 2024. [Llm-score: Unveiling the power of large language models in text-to-image synthesis evaluation](#). *Advances in Neural Information Processing Systems*, 36.
- Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. 2024. [Cheap and quick: Efficient vision-language instruction tuning for large language models](#). *Advances in Neural Information Processing Systems*, 36.

- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. ChatGPT as a factual inconsistency evaluator for text summarization. *arXiv preprint arXiv:2303.15621*.
- Yuanjie Lyu, Zhiyu Li, Simin Niu, Feiyu Xiong, Bo Tang, Wenjin Wang, Hao Wu, Huanyong Liu, Tong Xu, and Enhong Chen. 2024. Crud-rag: A comprehensive chinese benchmark for retrieval-augmented generation of large language models. *arXiv preprint arXiv:2401.17043*.
- Oscar Mañas, Benno Krojer, and Aishwarya Agrawal. 2024. Improving automatic vqa evaluation using large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4171–4179.
- Rui Mao, Guanyi Chen, Xulang Zhang, Frank Guerin, and Erik Cambria. 2023. Gpteval: A survey on assessments of chatgpt and gpt-4. *arXiv preprint arXiv:2308.12488*.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.
- Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. 2024. Chartinstruct: Instruction tuning for chart comprehension and reasoning. *arXiv preprint arXiv:2403.09028*.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209.
- Alexandre Matton, Tom Sherborne, Dennis Aumiller, Elena Tommasone, Milad Alizadeh, Jingyi He, Raymond Ma, Maxime Voisin, Ellen Gilsenan-McMahon, and Matthias Gallé. 2024. [On leakage of code generation evaluation datasets](#).
- Timothy R McIntosh, Teo Susnjak, Tong Liu, Paul Watters, and Malka N Halgamuge. 2024. Inadequacies of large language model benchmarks in the era of generative artificial intelligence. *arXiv preprint arXiv:2402.09880*.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. [Large language models: A survey](#).
- Abhika Mishra, Akari Asai, Vidhisha Balachandran, Yizhong Wang, Graham Neubig, Yulia Tsvetkov, and Hannaneh Hajishirzi. 2024. Fine-grained hallucination detection and editing for language models. *arXiv preprint arXiv:2401.06855*.
- Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig, and Yang You. 2024. [Mixeval: Deriving wisdom of the crowd from llm benchmark mixtures](#).
- OpenAI. 2023. [Gpt-4 technical report](#).
- Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B Hashimoto. 2023. Proving test set contamination in black box language models. *arXiv preprint arXiv:2310.17623*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Md Rizwan Parvez. 2024. Evidence to generate (e2g): A single-agent two-step prompting for context grounded and retrieval augmented reasoning. *arXiv preprint arXiv:2401.05787*.
- Md Rizwan Parvez, Wasi Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. [Retrieval augmented code generation and summarization](#). In *Findings*

- of the Association for Computational Linguistics: *EMNLP 2021*, pages 2719–2734, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Md Rizwan Parvez, Tolga Bolukbasi, Kai-Wei Chang, and Venkatesh Saligrama. 2019. [Robust text classifier on test-time budgets](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1167–1172, Hong Kong, China. Association for Computational Linguistics.
- Md Rizwan Parvez, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2018. [Building language models for text with named entities](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2373–2383, Melbourne, Australia. Association for Computational Linguistics.
- Md Rizwan Parvez and Kai-Wei Chang. 2021. [Evaluating the values of sources in transfer learning](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5084–5116, Online. Association for Computational Linguistics.
- Md Rizwan Parvez, Jianfeng Chi, Wasi Uddin Ahmad, Yuan Tian, and Kai-Wei Chang. 2023. [Retrieval enhanced data augmentation for question answering on privacy policies](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 201–210, Dubrovnik, Croatia. Association for Computational Linguistics.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. [Red teaming language models with language models](#). *arXiv preprint arXiv:2202.03286*.
- Sarah Masud Preum, Md Rizwan Parvez, Kai-Wei Chang, and John Stankovic. 2018. [A corpus of drug usage guidelines annotated with type of advice](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Xiao Pu, Mingqi Gao, and Xiaojun Wan. 2023. [Summarization is \(almost\) dead](#). *arXiv preprint arXiv:2309.09558*.
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. [Is ChatGPT a general-purpose natural language processing task solver?](#) *arXiv preprint arXiv:2302.06476*.
- Haoyi Qiu, Wenbo Hu, Zi-Yi Dou, and Nanyun Peng. 2024. [Valor-eval: Holistic coverage and faithfulness evaluation of large vision-language models](#). *arXiv preprint arXiv:2404.13874*.
- Qwen2. 2024. [Hello qwen2](#).
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. [Improving language understanding by generative pre-training](#).
- Raian Rahman, Rizvi Hasan, Abdullah Al Farhad, Md. Tahmid Rahman Laskar, Md. Hamjajul Ashmafee, and Abu Raihan Mostofa Kamal. 2023. [Chartsumm: A comprehensive benchmark for automatic chart summarization of long and short summaries](#). *Proceedings of the Canadian Conference on Artificial Intelligence*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don't know: Unanswerable questions for squad](#). *arXiv preprint arXiv:1806.03822*.
- Mathieu Ravaut, Bosheng Ding, Fangkai Jiao, Hailin Chen, Xingxuan Li, Ruochen Zhao, Chengwei Qin, Caiming Xiong, and Shafiq Joty. 2024. [How much are llms contaminated? a comprehensive survey and the llmsanitize library](#). *arXiv preprint arXiv:2404.00699*.
- Vipula Rawte, Amit Sheth, and Amitava Das. 2023. [A survey of hallucination in large foundation models](#). *arXiv preprint arXiv:2309.05922*.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. [Code llama: Open foundation models for code](#). *arXiv preprint arXiv:2308.12950*.

- Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. 2021. [How good is your tokenizer? on the monolingual performance of multilingual language models.](#) 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 3118–3135, Online. Association for Computational Linguistics.
- Mobashir Sadat, Zhengyu Zhou, Lukas Lange, Jun Araki, Arsalan Gundroo, Bingqing Wang, Rakesh Menon, Md Parvez, and Zhe Feng. 2023. [DelucionQA: Detecting hallucinations in domain-specific question answering.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 822–835, Singapore. Association for Computational Linguistics.
- Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure llm data contamination for each benchmark. *arXiv preprint arXiv:2310.18018*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- SambaNova. 2024. [Samba-coe v0.3: The power of routing ml models at scale.](#)
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. 2019. Socialliqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, Hyojung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson C. Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncarenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Miserlis Hoyle, and Philip Resnik. 2024. [The prompt report: A systematic survey of prompting techniques.](#)
- Sander Schulhoff, Jeremy Pinto, Anam Khan, Louis-François Bouchard, Chenglei Si, Svetlana Anati, Valen Tagliabue, Anson Kost, Christopher Carnahan, and Jordan Boyd-Graber. 2023. Ignore this title and hackaprompt: Exposing systemic vulnerabilities of llms through a global prompt hacking competition. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4945–4977.
- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.
- Shreya Shankar, JD Zamfirescu-Pereira, Björn Hartmann, Aditya G Parameswaran, and Ian Arawjo. 2024. Who validates the validators? aligning llm-assisted evaluation of llm outputs with human preferences. *arXiv preprint arXiv:2404.12272*.
- Chenhui Shen, Liying Cheng, Xuan-Phi Nguyen, Yang You, and Lidong Bing. 2023. Large language models are not yet human-level evaluators for abstractive summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4215–4233.
- Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang, Yifan Wang, Yujiu Yang, and Wai Lam. 2024. A thorough examination of decoding methods in the era of llms. *arXiv preprint arXiv:2402.06925*.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yongsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*.
- Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. Large language models are inconsistent and biased evaluators. *arXiv preprint arXiv:2405.01724*.
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang,

- Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*.
- Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. *arXiv preprint arXiv:1803.06643*.
- Yixuan Tang and Yi Yang. 2024. Multihoprag: Benchmarking retrieval-augmented generation for multi-hop queries. *arXiv preprint arXiv:2401.15391*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Simone Tedeschi, Felix Friedrich, Patrick Schramowski, Kristian Kersting, Roberto Navigli, Huu Nguyen, and Bo Li. 2024. Alert: A comprehensive benchmark for assessing large language models’ safety through red teaming. *arXiv preprint arXiv:2404.08676*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady Arkhangorodsky, Minjie Xu, Naomi White, and Patrick Lewis. 2024. Replacing judges with juries: Evaluating llm generations with a panel of diverse models. *arXiv preprint arXiv:2404.18796*.
- Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, et al. 2023a. On the robustness of ChatGPT: An adversarial and out-of-distribution perspective. *arXiv preprint arXiv:2302.12095*.
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023b. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*.
- Tong Wang, Ninad Kulkarni, and Yanjun Qi. 2024a. Less is more for improving automatic evaluation of factual consistency. *arXiv preprint arXiv:2404.06579*.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhua Chen. 2024b. [Mmlu-pro: A more robust and challenging multi-task language understanding benchmark](#).
- Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. 2023c. Learning to filter context for retrieval-augmented generation. *arXiv preprint arXiv:2311.08377*.
- Jason Wei, Yi Tay, and Quoc V Le. 2022a. Inverse scaling can become u-shaped. *arXiv preprint arXiv:2211.02011*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.
- Sheng-Lun Wei, Cheng-Kuang Wu, Hen-Hsen Huang, and Hsin-Hsi Chen. 2024. Unveiling selection biases: Exploring order and token sensitivity in large language models. *arXiv preprint arXiv:2406.03009*.
- Man-Fai Wong, Shangxin Guo, Ching-Nam Hang, Siu-Wai Ho, and Chee-Wei Tan. 2023. [Natural](#)

- language generation and understanding of big code for ai-assisted programming: A review. *Entropy*, 25(6):888.
- Honghan Wu, Minhong Wang, Jinge Wu, Farah Francis, Yun-Hsuan Chang, Alex Shavick, Hang Dong, Michael TC Poon, Natalie Fitzpatrick, Adam P Levine, et al. 2022. A survey on clinical natural language processing in the united kingdom from 2007 to 2022. *NPJ digital medicine*, 5(1):186.
- Minghao Wu and Alham Fikri Aji. 2023. Style over substance: Evaluation biases for large language models. *arXiv preprint arXiv:2307.03025*.
- Congying Xia, Chen Xing, Jiangshu Du, Xinyi Yang, Yihao Feng, Ran Xu, Wenpeng Yin, and Caiming Xiong. 2024. Fofo: A benchmark to evaluate llms' format-following capability. *arXiv preprint arXiv:2402.18667*.
- Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. *arXiv preprint arXiv:2402.13178*.
- Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. 2024. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. [Qwen2 technical report](#).
- Linyi Yang, Shuibai Zhang, Libo Qin, Yafu Li, Yidong Wang, Hanmeng Liu, Jindong Wang, Xing Xie, and Yue Zhang. 2022. Glue-x: Evaluating natural language understanding models from an out-of-distribution generalization perspective. *arXiv preprint arXiv:2211.08073*.
- Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao Gong, Yang Shen, et al. 2023a. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models. *arXiv preprint arXiv:2303.10420*.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023b. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.
- Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2023. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.
- Yuheng Zha, Yichi Yang, Ruichen Li, and Zhit-ing Hu. 2023. Alignscore: Evaluating factual consistency with a unified alignment function. *arXiv preprint arXiv:2305.16739*.
- Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja, Dylan Slack, Qin Lyu, et al. 2024a. A careful examination of large language model performance on grade school arithmetic. *arXiv preprint arXiv:2405.00332*.

- Pan Zhang, Xiaoyi Dong Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, et al. 2023. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.
- Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024b. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics*, 12:39–57.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023a. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. 2023b. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023a. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.
- Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong Wen, and Jiawei Han. 2023b. Don't make your llm an evaluation benchmark cheater. *arXiv preprint arXiv:2311.01964*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023a. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.
- Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. 2024. *Dyval 2: Dynamic evaluation of large language models by meta probing agents*.
- Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Neil Zhenqiang Gong, Yue Zhang, et al. 2023b. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv preprint arXiv:2306.04528*.
- Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou, and Ji-Rong Wen. 2023c. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107*.
- Ziyu Zhuang, Qiguang Chen, Longxuan Ma, Mingda Li, Yi Han, Yushan Qian, Haopeng Bai, Zixian Feng, Weinan Zhang, and Ting Liu. 2023. Through the lens of core competency: Survey on evaluation of large language models. *arXiv preprint arXiv:2308.07902*.
- Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. 2023. Red teaming ChatGPT via jailbreaking: Bias, robustness, reliability and toxicity. *arXiv preprint arXiv:2301.12867*.
- Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can large language models transform computational social science? *Computational Linguistics*, pages 1–55.

A Appendix

A.1 Benchmarking Datasets

General Capability Benchmarks: To benchmark the performance of LLMs, researchers typically use a set of widely recognized datasets. These common benchmarks are employed by authors upon the release of an LLM to evaluate its general capabilities. One of the most frequently used benchmarks is the MMLU benchmark (Hendrycks et al., 2020b), which assesses LLMs' overall knowledge and reasoning abilities

Criteria	Challenges & Limitations	Description
Reproducibility	Missing Experimental Details	Lack of documentation on the data subsets used for evaluation, which few-shot examples added to the prompt, what decoding parameters are used, etc., will impact reproducibility.
	Not Releasing the Data	The detailed prompt as well as the response generated by the LLMs are often missing.
	Code Unavailable	Many studies do not release the necessary codes (e.g., parsing scripts). This may impact reproducibility of the results.
	Model Updates and Depreciation	Continuous updates to the closed-source models (alongside possible depreciation of the models) will create challenges for reproducing previous results.
Reliability	Not Documenting Model Versions	The exact version of the model being used is often missing. This creates another reproducibility concern.
	Data Integrity	Incorrect gold labels and outdated benchmark datasets compromise evaluation reliability.
	Unfair Comparisons	Comparing models evaluated on the full dataset against the subset of a dataset, different few-shot examples being selected, etc.
	Contamination	LLMs may encounter evaluation data during pre-training, leading to contamination.
	Prompt Hacking	Manipulating input prompts to elicit desired responses can undermine fair evaluation.
	Transparency in Decoding Parameters	Lack of transparency in how decoding parameters are selected can lead to unfair comparisons.
Robustness	Evaluation Methodology and Metrics	Reliance on string-based metrics and automated evaluation methods without proper validation can lead to unreliable results.
	Limiting Evaluation to Certain Benchmarks	Evaluating LLMs only on a set of common benchmarks does not ensure generalizability.
	Lack of Diversity in Prompts and Parameters	Most existing research used only a single prompt while also not tuning any of the decoding parameters, restricting the robustness of the evaluation.
	Insufficient Evaluation Metrics	Lack of correlation between existing evaluation metrics impacts evaluation robustness.

Table 5: Challenges and Limitations in terms of Reproducibility, Reliability, and Robustness in LLM Evaluation.

across various subjects. Other common benchmarks focus primarily on evaluating the common-sense reasoning capabilities of LLMs (Wei et al., 2022a), such as HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020), SIQA, (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021), OpenBookQA (Mihaylov et al., 2018), ARC (Clark et al., 2018). In addition, the TruthfulQA dataset (Lin et al., 2021) is used to measure the truthfulness of an LLM, while the TyDi QA dataset (Clark et al., 2020) is used for evaluating the information seeking question answering capability across diverse languages. For assessing coding capabilities, the HumanEval (Chen et al., 2021) and the MBPP (Austin et al., 2021) are two widely used benchmarks. Additional problem-solving datasets include APPS (Hendrycks et al., 2021), CodeContests (Li et al., 2022), and xCodeEval (Khan et al., 2023), among others.

Specialized Benchmarks: There are also specialized benchmarks that measure specific capabilities of LLMs. For instance, the MT-Bench (Zheng et al., 2024) evaluates whether LLMs can properly engage in conversations, the RewardBench (Lambert et al., 2024) assesses the performance of reward models. Other specialized benchmarks like the AlpacaEval¹⁰ evaluates the instruction following capabilities (Zhou et al., 2023a) of LLMs, the Open Medical-LLM Leaderboard¹¹ evaluates the biomedical capabilities of LLMs, HHEM¹² leaderboard for hallucination detection (Mishra et al., 2024; Sadat et al., 2023), BigCodeBench¹³

¹⁰https://tatsu-lab.github.io/alpaca_eval/

¹¹<https://huggingface.co/blog/leaderboard-medicalllm>

¹²<https://huggingface.co/spaces/vectara/leaderboard>

¹³<https://huggingface.co/blog/leaderboard-bigcodebench>

and LiveCodeBench¹⁴ for code generation capability evaluation, SWE-bench (Jimenez et al., 2023) for software engineering capability evaluation. The recently proposed FOFO benchmark Xia et al. (2024) measures language models' ability to adhere to the requested formats in prompts across different domains. Moreover, there are also some specialized benchmarks that are used for LLM safety¹⁵ (Chao et al., 2024) and red teaming¹⁶ (Tedeschi et al., 2024) evaluation. The ability to understand both language and vision is another interesting capability of recently proposed LLMs (Bai et al., 2023; Chen et al., 2023a; Dai et al., 2024; Liu et al., 2023b, 2024a; Luo et al., 2024; Ye et al., 2023b; Zhang et al., 2023; Zhu et al., 2023a). This has led to the development of many multi-modal benchmarks (Chen et al., 2024b; Fu et al., 2023a, 2024a; Guan et al., 2023; Li et al., 2023a,b,d; Liu et al., 2024a, 2023e; Lu et al., 2022; Qiu et al., 2024; Yu et al., 2023). These benchmarks study the multimodal capabilities of LLMs across various domains, such as math and reasoning (Lu et al., 2023; Yue et al., 2023), science diagrams (Kembhavi et al., 2016), chart understanding and reasoning (Islam et al., 2024b; Masry et al., 2022, 2024; Rahman et al., 2023), document understanding (Mathew et al., 2021).

Other Diverse Benchmarks: To enable a more comprehensive evaluation of LLMs across a wide range of scenarios, some studies also focused on introducing new benchmarks covering various aspects, such as HELM (Liang et al., 2022), PromptBench (Zhu et al., 2023b), OpenLLM¹⁷, MixEval (Ni et al., 2024), etc. These benchmarks cover diverse tasks and usually include existing benchmarking datasets (e.g., MMLU, HellaSwag, BoolQ (Clark et al., 2019), etc.). Additionally, despite the availability of numerous benchmarks (both general and specialized), existing widely-used benchmarks still do not cover the full variety of tasks (Parvez et al., 2018; Preum et al., 2018). Therefore, some researchers have independently evaluated LLMs using additional di-

¹⁴<https://huggingface.co/blog/leaderboard-livecodebench>

¹⁵<https://huggingface.co/spaces/AI-Secure/llm-trustworthy-leaderboard>

¹⁶<https://huggingface.co/spaces/HaizeLabs/red-teaming-resistance-benchmark>

¹⁷https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard

verse NLP datasets and tasks (Bang et al., 2023; Kocoń et al., 2023; Laskar et al., 2023a; Qin et al., 2023). They also employed domain-specific benchmarks in fields such as biomedicine (Jahan et al., 2023, 2024), finance (Guo et al., 2023a; Li et al., 2023e), language-specific (Abdelali et al., 2024; Ahuja et al., 2023; Kabir et al., 2023; Khondaker et al., 2023; Lai et al., 2023; Liu et al., 2023c), social science (Ziems et al., 2024), coding (Liu et al., 2024b), and information retrieval (Zhu et al., 2023c). In addition to that, ethics, bias, toxicity, robustness, and trustworthiness are also independently evaluated by researchers across various datasets (Hendrycks et al., 2020a; Liu et al., 2023a; McIntosh et al., 2024; Rawte et al., 2023; Sun et al., 2024; Wang et al., 2023a; Yang et al., 2022; Zhuo et al., 2023).

A.2 Prominent LLMs

The impressive success of ChatGPT has led to the development of many LLMs in recent years. Since there are hundreds of LLMs being released in recent years (Zhao et al., 2023a), we only discuss some of the prominent LLMs that achieved top rankings in various public leaderboards recently. LLMs can be categorized into two parts: *Closed-Source LLMs*: only available for use through the API or web interface, and (ii) *Open-Source LLMs*: where the pre-trained weights of the model are available that allow further training of such models. Below, we present some prominent LLMs in these two categories.

A.2.1 Closed Source LLMs

In the following, we categorize LLMs based on the organizations that develop these LLMs:

OpenAI models (OpenAI, 2023):

- **GPT-3.5:** This model is an iteration of the GPT-3 architecture, emphasizing improvements in response quality through the application of the reinforcement learning from human feedback (RLHF) technique. GPT-3.5 is known for its robust performance in zero-shot tasks, where no specific training examples are provided during the task execution. This model has been instrumental due to its strong foundational capabilities in understanding and generating human-like text.
- **GPT-4:** It extends GPT-3.5's capabilities by

incorporating multimodal functionalities, allowing the model to process not just text but also visual inputs. This advancement significantly broadens its applicational scope, making it adept at handling more complex tasks that require an understanding of both textual and visual information. It features enhanced safety protocols and a sophisticated training regime that includes a safety reward signal during its reinforcement learning phase.

- **GPT-4 Turbo:** This version builds upon GPT-4's foundation with substantial upgrades in computational efficiency and functionality. GPT-4 Turbo boasts an increased model capacity and an extended knowledge base that encompasses more recent data up to April 2023. It features a longer context window of up to 128,000 tokens and includes significant improvements in the model's economy and output consistency.
- **GPT-4o:** OpenAI's most sophisticated model, GPT-4o ("o" for "omni") is a multimodal powerhouse capable of handling both text and image inputs to generate text outputs. It improves upon GPT-4 Turbo by offering double the text generation speed and reducing operational costs by 50%.

Google models:

- **PaLM-2:** Released by Google in 2023, it is an advanced large language model that builds on the foundations set by its predecessor, the original PaLM. This iteration incorporates a sophisticated 'mixture of objectives' technique, allowing it to surpass the capabilities of the earlier model significantly (Anil et al., 2023).
- **Gemini:** It is a multimodal model developed by Google in December 2023, to understand and process a variety of information types, including text, images, audio, and video, seamlessly. Gemini's architecture allows it to perform exceptionally across multiple platforms, from large-scale data centers to mobile devices, adapting efficiently to the needs of different applications. This model sets new benchmarks in AI with its ability to excel in tasks that require complex multimodal integrations (Team et al., 2023).

Anthropic Models: *The Claude series* models, developed by Anthropic, represent a series of advanced language models designed to enhance user interaction through natural language understanding and generation. Starting with the original Claude, which excelled in tasks like summarization and creative writing, each subsequent model—Claude Instant, Claude 2.0, and the Claude 3 family (Haiku, Sonnet, and Opus)—has introduced significant improvements in processing speed, reasoning capabilities, and multimodal functionality. These models have a variety of uses, from quick response generation in Claude Instant to sophisticated multimodal understanding in Claude 3 Opus, showcasing their versatility and advanced AI technology to meet different user and enterprise needs¹⁸. The latest model in the Claude-3 series is the Claude-3.5-Sonnet¹⁹ model.

A.2.2 Open Source LLMs

We similarly categorize the open-source LLMs based on the organizations that develop them:

Meta Models:

- **Llama:** Launched in February 2023 by Meta AI, Llama was the first in the Llama series, showcasing strong performance on a range of natural language processing tasks. It competed well against larger models like GPT-3 with a smaller parameter size and was made available under a non-commercial license, primarily for academic research (Touvron et al., 2023a).
- **Llama 2:** Released in July 2023, Llama 2 improved on its predecessor by expanding model sizes up to 70 billion parameters. It maintained the original architecture but included better training data and enhanced functionality. Notably, Llama 2 was more accessible, available for both academic and some commercial uses (Touvron et al., 2023b).
- **Llama 3:** In April 2024, Meta AI introduced Llama 3²⁰, the most advanced version with up to 70 billion parameters. This version added longer context capabilities and

¹⁸<https://www.anthropic.com/news/claude-3-family>

¹⁹<https://www.anthropic.com/news/claude-3-5-sonnet>

²⁰<https://llama.meta.com/llama3/>

improved multimodal functions, marking a significant advancement in AI technology application across various fields.

Mistral Models: Mistral AI, founded in April 2023, is specialized in the development of open-source large language models. Rapidly gaining recognition in the AI industry, Mistral AI emphasizes the importance of open-source software, providing a viable alternative to proprietary models. The company has released several models, including Mistral 7B, Mixtral 8x7B, and Mixtral 8x22B, which are known for their high performance and innovation in the use of mixture of experts architectures (Cai et al., 2024; Jiang et al., 2023). Codestral 22B, introduced on May 29, 2024, is a pioneering code generation model designed to enhance coding efficiency across more than 80 programming languages. With its specialized focus and lightweight architecture, Codestral significantly outperforms other leading models on the HumanEval FIM benchmark, making it a critical tool for developers seeking advanced AI-assisted coding capabilities.

Alibaba Models: QWEN series models are transformer-based large language models developed by Alibaba Cloud (Bai et al., 2023). These models, pre-trained on diverse data sources including web texts, books, code, and more, come in various sizes ranging from 0.5 billion to 110 billion parameters. Qwen models support long context lengths and demonstrate strong performance on multiple Chinese and English evaluation tasks, including common-sense reasoning, code, and mathematics. The latest versions, Qwen 1.5 and Qwen 2, offer significant improvements in chat model performance, multilingual support, and stable support for up to 32K context length. With a comprehensive vocabulary of over 150K tokens, Qwen models are designed to handle multiple languages effectively, making them a versatile tool for various AI applications.

Microsoft Models: The Phi series (Abdin et al., 2024) by Microsoft consists of small language models (SLMs) designed to provide high performance with lower computational requirements. The newly announced Phi-3 family includes models like Phi-3-mini, Phi-3-small, and Phi-3-medium, ranging from 3.8 billion to 14 billion parameters. These models excel in various benchmarks, offering capabilities similar to larger mod-

els but in a smaller, more cost-effective package. Phi-3 models are particularly suited for simpler tasks, local device operations, and environments with limited resources, making AI more accessible and efficient for diverse applications. They are available through Microsoft Azure AI Model Catalog, Hugging Face, and as NVIDIA NIM microservices. Several followup works extends Phi-models or their synthetic data into multilingual space such as (Boughorbel et al., 2024).

Technology Innovation Institute Models: Technology Innovation Institute release the Falcon series models (Almazrouei et al., 2023), such as the Falcon 2 series that include models with parameter sizes of 1.3B, 7.5B, 40B, and 180B. These models are notable for their use of the REFINEDWEB dataset. Falcon models are designed for both research and commercial use, with Falcon 2 models featuring multilingual and multimodal capabilities, including vision-to-language. The Falcon 180B model, in particular, is accessible under a royalty-free license.

Cohere Models: Cohere offers a variety of advanced large language models designed for multiple use cases, including text generation, embeddings, and reranking. The Command family models, such as Command R+ and Command R, excel in conversational tasks and complex workflows like code generation and retrieval-augmented generation (RAG)²¹ (Alonso et al., 2024; Chen et al., 2024a; Gao et al., 2023b; Lewis et al., 2020; Liu et al., 2023d; Lyu et al., 2024; Parvez et al., 2021, 2023; Tang and Yang, 2024; Wang et al., 2023c; Xiong et al., 2024). The Embed models enhance search, classification, and clustering capabilities with both English and multilingual support. The Rerank models improve search algorithms by reorganizing results based on specified parameters. Cohere models are accessible across platforms like Amazon SageMaker, Microsoft Azure, and Oracle GenAI Service, enabling seamless integration into diverse applications and retrieval augmented generation.

Google Gemma Models: While early LLMs released by Google's are mostly closed-source (e.g., PaLM-2, Gemini, etc.), Google has also recently released some lightweight open-source LLMs,

²¹<https://cohere.com/command>

named as Gemma²² family LLMs, that also have multimodal capabilities²³.

A.3 Prompting Techniques

Prompts can be designed in various ways (Brown et al., 2020; Chung et al., 2022; Islam et al., 2024a; Parvez, 2024; Schulhoff et al., 2024; Wei et al., 2022b), as stated below:

- **In-Context Learning (Zero-shot):** It means that the prompt used to interact with the model contains no examples or demonstrations. The model relies on its pre-existing knowledge, obtained from its initial training on diverse data, to generate a response or perform the task based solely on the instructions given. For example, “classify the sentence as biased or unbiased text”.
- **In-Context Learning (Few-shot):** It means that the prompt used to interact with the model includes a small number of examples or demonstrations. The model uses these examples to quickly adapt and understand how to perform a specific task, leveraging the details within these examples. This technique allows the model to extend its pre-existing knowledge to new tasks by closely analyzing the limited examples given. For instance, classify the sentence as biased or unbiased based on a few similar examples provided.
- **Chain-of-Thought Prompting (CoT):** This technique encourages models to generate intermediate reasoning steps before arriving at a final answer, mimicking a human-like problem-solving approach. This can be combined with few-shot prompting to achieve better results on more complex tasks. For example, if asked to determine whether the number "15" is odd or even, the model might outline its reasoning as follows: "An even number is divisible by 2 without a remainder. 15 divided by 2 is 7 with a remainder of 1. Therefore, 15 is an odd number." This step-by-step explanation helps clarify the model's thought process and supports its conclusion.
- **Decomposition Techniques:** These techniques break down complex problems into

simpler sub-problems that can be solved sequentially by the GenAI model. Each component of the problem is addressed individually, and the solutions are integrated to form a comprehensive response. Decomposition is especially useful in tasks that require layered reasoning or have multiple steps. For example, in solving a math word problem, decomposition might involve separately calculating the distances each person travels and then combining these calculations to determine when they meet.

- **Role-based and Style-based Prompting:** In these techniques prompts are designed to induce a specific style or persona in the model's responses. By specifying a role (e.g., a scientist explaining a concept) or a style (e.g., formal or poetic), users can guide the tone and formality of the AI's output. This technique is valuable in applications requiring genre-specific content generation or when the output needs to fit a particular communicative context.
- **Prompt chaining:** It is a technique where a complex task is divided into simpler sub-tasks, each addressed by its own prompt. The response from one prompt is used as the input for the next, creating a sequential chain of prompts that gradually build towards the final answer. This method enhances the performance and reliability of large language models by breaking down tasks into manageable parts, making it easier to control and refine the model's responses at each step. For example, in a document analysis task, the first prompt might extract key facts from a text, and the second prompt would use these facts to generate a summary.
- **Tree of Thoughts (ToT):** It is a technique that structures problem-solving into a tree of possible solutions. It uses strategies like breadth-first or depth-first search to evaluate each potential solution path. For example, in solving a puzzle, ToT might explore different moves to find the quickest solution path.
- **Directional Stimulus Prompting (DSP) :** It is a technique that enhances how large language models (LLMs) respond to tasks by using dynamically generated prompts. A

²²<https://storage.googleapis.com/deepmind-media/gemma/gemma-report.pdf>

²³<https://huggingface.co/blog/paligemma>

secondary, tuneable model creates specific hints that guide the main, unchangeable LLM to produce more targeted and relevant outputs. This method uses reinforcement learning to refine these prompts based on how well they perform, making DSP a more adaptive and precise approach compared to standard prompting techniques. For instance, in summarizing complex documents, DSP might generate a prompt like "Summarize focusing on economic impacts," guiding the LLM to tailor its output specifically to the economic aspects mentioned in the text.

- **Multimodal Prompting:** Extending beyond text, multimodal prompting involves using inputs like images, audio, or video along with textual descriptions. This technique leverages the model's capability to process and integrate information from diverse data types, enhancing its applicability in scenarios where multiple forms of data are available. For example, interpret a scene from a video by analyzing both the spoken dialogue and the visual content to determine the mood of the conversation.
- **Meta-Prompting:** It involves creating prompts that instruct the AI to generate or refine its prompts, essentially using AI to improve the efficiency and effectiveness of prompt engineering. This recursive use of prompting can lead to more dynamic and contextually adaptive AI behaviors. For example, ask the AI to optimize a prompt that instructs another AI to summarize news articles, thereby refining the instructions to enhance summary relevance and conciseness.

A.4 Decoding Parameters

There are various decoding parameters that are required to be set. For instance:

- **Temperature:** It is used to control the randomness of the output. It is typically between 0 and 1. Lower values (e.g., 0.1) make the model more deterministic and focused on the most likely next token, while higher values (e.g., 0.9) introduce more randomness and diversity.
- **Beam Size:** It refers to the number of beams in Beam Search (Freitag and Al-Onaizan,

2017), a decoding strategy that keeps track of multiple possible sequences (beams) at each step of generation to find the most likely sequence. A higher number of beams usually leads to more accurate results but at the cost of increased computation.

- **Top-K:** The number of top probable tokens to consider. For example, if $K=10$, the model will choose the next token only from the top 10 most likely tokens.
- **Top-P:** The cumulative probability threshold. For example, if $P=0.9$, the model will sample from the smallest set of tokens whose combined probability is at least 90%.
- **Maximum Output Tokens:** It sets the maximum number of tokens to generate.

A.5 Parsing Script Design

While there are various evaluation software (Biderman et al., 2024; Dalvi et al., 2023) currently available, they are limited to certain scenarios (e.g., limited to certain datasets and benchmarks, prompts, etc.). Thus, for the evaluation of LLMs across diverse settings, researchers often require to write parsing scripts. We present some scenarios in Table 6 to demonstrate why parsing script is required for such cases and the importance of validating parsing scripts.

A.6 Evaluation Approach

A.6.1 Automatic Evaluation

To provide a high-level overview, automatic evaluation for LLMs can be divided into the following:

Language Modeling: Perplexity (Jelinek et al., 1977) is widely used to study the performance of auto-regressive language models. It measures how confidently a model predicts the next word in a sequence, with the assumption that lower perplexity indicates better performance. Hence, perplexity has been historically used to assess the language model's capability to generate a coherent language and is also useful to quickly compare different models or checkpoints.

Discriminative Tasks: For tasks involving class prediction, post-processing using a parsing script is usually required to extract answers from the LLM-generated responses to compare against gold labels. In this context, metrics such as Exact Match, Accuracy, Precision, Recall, and F1,

Scenario 1: For the response generated, designing a parsing script to extract the answer “Lionel Messi” is straightforward. However, the parsing script should also be robust to cover cases like abbreviations, uppercase-lowercase sensitivity, punctuations, synonyms, stemming, lemmatization, paraphrases, etc.

Prompt: Which player has won the best player award in Fifa world cup 2022?

Sample LLM Response (GPT 4o): Lionel Messi won the Best Player award (Golden Ball) in the FIFA World Cup 2022. He was instrumental in leading Argentina to victory in the tournament, culminating in their triumph in the final against France.

Correct Answer: Lionel Messi

Scenario 2: While Extraction of the answer “Lionel Messi” is required, due to the LLM knowledge-cut-off date of September 2021, it may answer about 2018. However, the target answer “Lionel Messi” is also in the output and so if the parsing script only parses the target answer then it may consider the response as correct whereas the response is wrong.

Prompt: Which player has won the best player award in the last Fifa world cup?

Sample LLM Response (Older ChatGPT 3.5 having knowledge cut-off date of September 2021): The Best Player award (Golden Ball) in the previous FIFA World Cup, which was held in 2018 in Russia, was won by Luka Modric from Croatia. Prior to that, Lionel Messi had won it in 2014.

Correct Answer: Lionel Messi

Table 6: Some examples of LLM-generated response requiring parsing script to extract the target answer. For Scenario 2, human evaluation is usually needed to ensure accurate parsing of the answer.

are usually utilized in discriminative tasks (Bang et al., 2023; Laskar et al., 2023a; Qin et al., 2023). Since metrics like exact match have several limitations (e.g., they do not consider the synonym of the gold label), various metrics for certain tasks (e.g., question answering (Bulian et al., 2022; Chen et al., 2020; Li et al., 2024; Mañas et al., 2024)) are proposed.

Generative Tasks: For generative tasks such as summarization or machine translation, parsing scripts are usually not required (Jahan et al., 2024; Laskar et al., 2023a) and so the full response generated by LLMs are compared against the gold reference. In this regard, *ROUGE* (Lin, 2004) and *BLEU* (Papineni et al., 2002) which are based on n-gram word matching are widely used. Meanwhile, various contextualized similarity (Laskar et al., 2020; Parvez and Chang, 2021) metrics (e.g., *BERTScore* (Zhang et al., 2019), *BARTScore* (Yuan et al., 2021), *AlignScore* (Wang et al., 2024a; Zha et al., 2023)) are also utilized that do not depend on word-based similarity measures.

A.6.2 Human Evaluation

Since LLMs generate human-like responses, it is often required to conduct qualitative evaluation of their responses. Earlier, qualitative evaluation of model-generated responses in terms of fluency, coherence, and informativeness were very popular

(Laskar et al., 2022b). However, with LLMs usually generating informative, fluent, and coherent response (Bang et al., 2023; Kocoń et al., 2023; Laskar et al., 2023a; Qin et al., 2023), the evaluation of factual consistency of LLM-generated responses has become more important recently (Fu et al., 2023b). Moreover, qualitative evaluation to compare between LLM-generated responses via leveraging humans based on the Elo rating system (Zheng et al., 2024) has gained a lot of attention.

Elo Rating: Elo rating works by comparing LLMs in pairwise “A vs B” comparisons, where each model is assigned an initial numerical rating (Boubdir et al., 2023; Zhao et al., 2023b). The outcome of each comparison adjusts these ratings based on the Elo algorithm: if a model performs better than expected, its rating increases; if it performs worse, its rating decreases. The expectation of a model’s performance is calculated using its rating relative to its opponent’s, adjusted by a factor that represents the sensitivity of expected scores to differences in ratings. To ensure a robust evaluation of LLMs using the Elo benchmark, it’s important to follow key indicators like reliability and transitivity (Boubdir et al., 2023). Reliability keeps Elo ratings consistent across various comparison sequences and prevents them from being overly sensitive to changes in hyperparameters, such as the K-factor. Transitivity is cru-

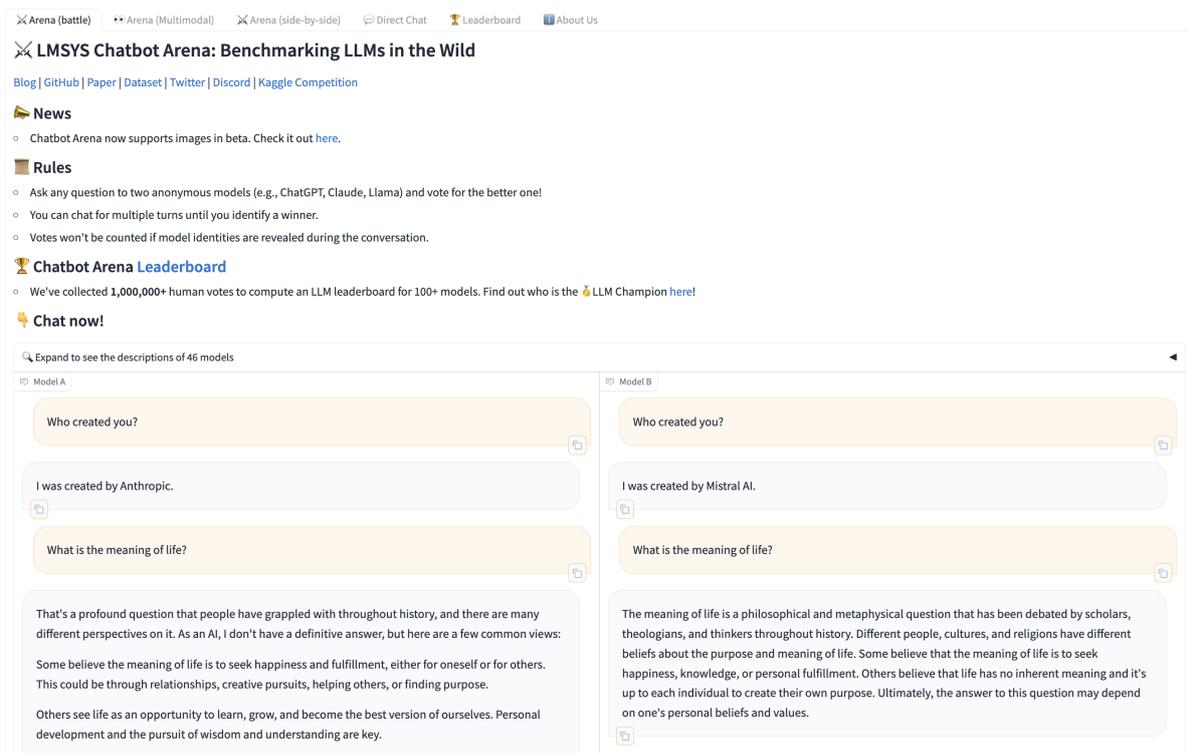


Figure 5: Ownership attack for blind evaluation on LLMs: Reviewers can pose any ownership-related questions and select their preferred model solely based on the ownership of the model. LMSys doesn't count votes if the model's identities are revealed during conversation

cial, indicating that if model A is rated higher than model B, and model B is rated higher than model C, model A should logically rank above model C. Extensive testing with both synthetic and real-world data is essential to verify that Elo scores accurately and stably reflect model performance (Boubdir et al., 2023). This involves making precise adjustments to the comparison order, selecting hyperparameters carefully, and utilizing numerous permutations to ensure outcome consistency. Due to the sensitive nature of the Elo rating system towards the order in which the updates were performed, Zheng et al. (2024) used the Bradley-Terry (BTL) model for their chatbot arena ranking. It is observed that model A can have a higher win rate than model B both empirically and statistically but a lower Elo rating. Since win rate serves as the stand-in measure for the probability of a model being better than another, this signifies the findings by Boubdir et al. (2023) that Elo rating is non-transitive with or without (BTL). On the other hand, BTL-based rating is tolerant to an imbalanced number of votes per model as shown by (Zheng et al., 2024), they also propose a different probability of win rates that are derived from

the ratings found from BTL which is transitive but doesn't correlate with the empirical win rates.

Elo hacking: Crowdsourced Elo-based ranking has gained popularity through the LMSys leaderboard²⁴ and has been accepted by various organizations, prompting them to release their LLMs early into this ecosystem for human evaluation. However, such setups can be easily exploited on a large scale using simple techniques. Figure 5 illustrates how someone can initially bypass the blind scoring mechanism through ownership hacking. Additionally, the evaluation of knowledge bases is not easily tracked, making votes on highly complex reasoning questions equivalent to those on simpler queries. Furthermore, upon the release of a popular model, systematic attacks or boosting can be initiated through ownership hacking. In addition to that, considering same score for *tie* and *both-bad* can significantly change leaderboard position. We recommend to use *tie* as 0.5 point and *both-bad* as 0 point.

²⁴<https://huggingface.co/spaces/lmsys/chatbot-arena-leaderboard>

A.6.3 LLMs as Evaluators

Since human evaluation is time-consuming (Laskar et al., 2023c,d) and difficult to reproduce, the instruction-following capabilities of LLMs have also inspired researchers to use certain LLMs as a judge to evaluate the responses generated by other LLMs (Chern et al., 2024; Fu et al., 2023b; Gao et al., 2023a; Hada et al., 2023; Huang et al., 2024a; Kenton et al., 2024; Kim et al., 2024b; Kobayashi et al., 2024; Kocmi and Federmann, 2023; Lu et al., 2024; Luo et al., 2023; Perez et al., 2022; Shankar et al., 2024). While prior work mostly utilized general-purpose closed-source LLMs-as-a-judge, the recently proposed Prometheus 2 (Kim et al., 2024a) model is an open-source variant which is specifically trained for qualitative evaluation of model-generated responses and demonstrated higher correlation with humans.

However, research by (Wang et al., 2023b) and (Shen et al., 2023) has highlighted potential limitations in using LLM as evaluators, suggesting that while LLMs can excel in specific areas like translation quality and grammatical error correction (Kobayashi et al., 2024; Kocmi and Federmann, 2023), their effectiveness as evaluators may vary significantly across different tasks. Moreover, using closed-source LLMs as evaluators also have associated cost. This highlights the ongoing debate and research into the capabilities and limitations of LLMs as evaluators in diverse linguistic domains. Therefore, to use LLMs as evaluators, it is important to consider the following:

- **Consistency:** Ensuring consistent combinations of LLMs are used as evaluators when LLMs are used as juries to ensure consistency and reproducibility in assessments.
- **Bias and Hallucination Detection:** Developing methods to identify and mitigate bias and hallucinations in the outputs of LLM judges/juries to ensure the reliability and robustness of the evaluation.
- **Interpretability:** Enhancing the interpretability of LLM outputs (e.g., asking LLMs to provide reasoning/explanations) to improve understanding and trustworthiness of the evaluation.
- **Cost Efficiency:** Advancing the development of efficient LLMs to reduce costs.