

LLMs instead of Human Judges?

A Large Scale Empirical Study across 20 NLP Evaluation Tasks

Anna Bavaresco¹, Raffaella Bernardi², Leonardo Bertolazzi², Desmond Elliott³,
Raquel Fernández¹, Albert Gatt⁴, Esam Ghaleb¹, Mario Giulianelli⁵,
Michael Hanna¹, Alexander Koller⁶, André F. T. Martins⁷, Philipp Mondorf⁸,
Vera Neplenbroek¹, Sandro Pezzelle¹, Barbara Plank⁸, David Schlangen⁹,
Alessandro Suglia¹⁰, Aditya K Surikuchi¹, Ece Takmaz⁴, Alberto Testoni¹

¹University of Amsterdam, ²University of Trento, ³University of Copenhagen,
⁴Utrecht University, ⁵ETH Zürich, ⁶Saarland University, ⁷Universidade de Lisboa & Unbabel,
⁸LMU Munich & MCML, ⁹University of Potsdam, ¹⁰Heriot-Watt University,

Abstract

There is an increasing trend towards evaluating NLP models with LLM-generated judgments instead of human judgments. In the absence of a comparison against human data, this raises concerns about the validity of these evaluations; in case they are conducted with proprietary models, this also raises concerns over reproducibility. We provide JUDGE-BENCH, a collection of 20 NLP datasets with human annotations, and comprehensively evaluate 11 current LLMs, covering both open-weight and proprietary models, for their ability to replicate the annotations. Our evaluations show that each LLM exhibits a large variance across datasets in its correlation to human judgments. We conclude that LLMs are not yet ready to systematically replace human judges in NLP.

1 Introduction

For many natural language processing (NLP) tasks, the most informative evaluation is to ask humans to judge the model output. Such judgments are traditionally collected in lab experiments or through crowdsourcing, with either expert or non-expert annotators, as illustrated in Figure 1. Recently, there has been a trend towards replacing human judgments with automatic assessments obtained via large language models (LLMs) (Chiang and Lee, 2023; Wang et al., 2023; Liu et al., 2023, *inter alia*). In this setting, an LLM is prompted with the instructions of the evaluation task, and asked to generate a value from the space of possible annotations that suits the task. For example, the LLM could be instructed to rate a piece of text for perceived plausibility in a dialogue system response, on a scale from 1 to 5. This drastically reduces the evaluation effort and is claimed to yield more reliable results across multiple evaluation rounds by minimising annotator noise (Landwehr et al., 2023; Jiang et al., 2023b; Reiter, 2024; Dubois et al., 2024).

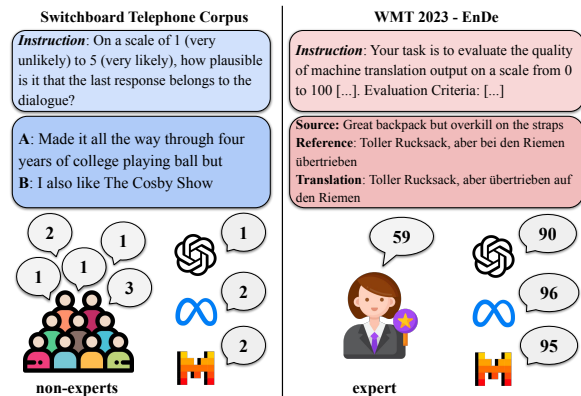


Figure 1: Evaluation by expert and non-expert human annotators and by LLMs for two tasks involving human-generated (left) and machine-generated text (right).

At the same time, the use of LLMs as judges of linguistic output raises new concerns. While LLMs have improved by leaps and bounds, they are *not* actual humans and may be prone to errors or systematic biases that differ from those of humans, especially on the increasingly subtle tasks that the NLP community has recently tackled, such as evaluating toxicity, or reasoning. This may distort evaluation results and lead to incorrect conclusions. The problem is aggravated by the fact that most LLMs do not disclose their training data, which makes it impossible to check for definitive data leakage from existing benchmarks, and undermines the ability to make broad, generalisable claims beyond the single specific dataset under analysis. Another, more implicit, type of data leakage can result from users unknowingly feeding training data into the system (Balloccu et al., 2024). Specifically for closed models such as OpenAI’s GPT series, there are thus serious reproducibility concerns, as LLMs may be retrained or retired at any time, to the point of rendering comparisons between this month’s and last month’s judgments invalid.

Previous studies offer mixed evidence regarding the reliability of LLMs compared to human annotators. Some research concludes that LLMs can generally serve as effective evaluators, correlating well with human judgments (Liu et al., 2023; Zheng et al., 2024; Chen et al., 2023; Verga et al., 2024; Törnberg, 2023; Huang et al., 2024; Naismith et al., 2023; Gilardi et al., 2023; Kocmi and Federmann, 2023b), albeit with some caveats (Wang et al., 2023; Wu and Aji, 2023; Hada et al., 2024; Pavlovic and Poesio, 2024). In some cases, LLM evaluators can also provide fine-grained evaluation beyond a single score, such as error spans (Fernandes et al., 2023; Kocmi and Federmann, 2023a). In contrast, some studies highlight significant flaws in LLM performance as evaluators (Koo et al., 2023; Zeng et al., 2024; Baris Schlicht et al., 2024), while other studies do not validate LLMs against human judgments (Jiang et al., 2023b; Landwehr et al., 2023). These discrepancies likely stem from the limitations of previous work, which typically relies on a few datasets and models, often restricted to closed-source proprietary models.

In this paper, we present JUDGE-BENCH to examine how well current LLMs can approximate human evaluators by analysing the correlation between model and human judgments on a large scale. We prompt 11 LLMs to generate judgments on 20 datasets with human annotations on a wide range of quality dimensions and tasks. Our evaluation goes beyond existing work by including a wide variety of *datasets* that differ in the type of task (e.g., translation, dialogue generation, etc.), the *property* being judged (e.g., coherence, fluency, etc.), the *type of judgments* (categorical or graded), and the *expertise of human annotators* (experts or non-experts). We evaluate the most recent open-weight and proprietary LLMs of different sizes.

We find that while some LLMs correlate well with human judgments on some datasets—indicating that they could indeed be used as valid surrogates—each tested LLM performs poorly on some others and exhibits significant variance across datasets. This means that current LLMs and/or their prompts *need to be calibrated against actual human judgments on every new dataset* to establish the validity of their evaluation scores. We observe a *decreasing gap between open and closed models*: the overall best-performing LLM in our evaluation is GPT-4o, with Llama3-70B coming in a close second; this seems promising with respect to the reproducibility of future evaluation efforts. We re-

lease our code base, which is intended as a living benchmark to ease extension and future evaluation, at the following link: <https://github.com/dmg-illc/JUDGE-BENCH>.

2 Construction of JUDGE-BENCH

In order to assess the capacity of LLMs to act as judges, we define JUDGE-BENCH, a comprehensive set of 20 datasets annotated by humans, for a range of quality dimensions. One key feature that differs across the datasets is the source of the data being evaluated as illustrated further in Figure 1, i.e., whether the items to be judged are generated by a model or produced by humans. In the former case (model-generated), the datasets concern tasks where we evaluate the performance of an NLP system; this includes classic downstream tasks such as machine translation, or dialogue response generation, as well as other less standard tasks for which automation has recently become an option thanks to LLMs, such as the generation of plans or logical arguments. In the latter case, the underlying goal is not to assess the performance of an NLP system, but rather properties of interest in human-generated language, such as grammaticality or toxicity.

Regardless of whether the instances to be judged are model- or human-generated, the judgments always concern a certain property of interest. The datasets we consider cover a wide span of properties, ranging from grammaticality and toxicity (as mentioned above) to coherence, factual consistency, and verbosity, *inter alia*. Many properties are relevant across multiple tasks (e.g., fluency and coherence), while others are more task-specific (e.g., success of a plan generator).

Our study focuses on English datasets or language pairs which include English as one of the languages. When sourcing the datasets, we keep track of whether the original annotation guidelines are available and whether the annotations are provided by expert or non-expert annotators. We retain all available individual annotations. Dataset information is summarised in Table 2, Appendix A. Furthermore, all 20 datasets are formatted following a precise data schema to facilitate the integration of additional datasets. We consider JUDGE-BENCH a living benchmark that could be easily extended.

3 Model Selection and Experiment Design

Models We select representative proprietary and open-weight models of various sizes that are widely

Type	Dataset (# properties judged)	GPT-4o	Llama-3-70B	Mixtral-8x22B	Gemini-1.5	Mixtral-8x7B	Comm-R+	σ	α
Categorical Annotations	cola (1)	0.34	0.56	0.54	0.44	0.55	0.12	0.16	-
	toxic-chat (2)	0.49 \pm 0.36	0.40 \pm 0.13	0.45 \pm 0.27	0.45 \pm 0.35	0.36 \pm 0.12	0.28 \pm 0.35	0.11	-
	llmbar-natural (1)	0.84	0.72	0.72	0.79	0.5	0.56	0.14	-
	llmbar-adversarial (1)	0.58	0.4	0.2	0.29	0.07	0.11	0.24	-
	cola-grammar (63)	0.47 \pm 0.22	0.29 \pm 0.24	0.28 \pm 0.23	0.26 \pm 0.24	0.21 \pm 0.18	0.13 \pm 0.14	0.13	-
	topical-chat (2)	0.05 \pm 0.07	-0.02 \pm 0.02	-0.03 \pm 0.04	-0.03 \pm 0.04	0.02 \pm 0.03	0.01 \pm 0.02	0.15	0.08
	roscoe-gsm8k (2)	0.59 \pm 0.35	0.65 \pm 0.27	0.62 \pm 0.38	0.6 \pm 0.24	0.58 \pm 0.36	0.0	0.23	-
	roscoe-esnli (2)	0.29 \pm 0.06	0.31 \pm 0.15	0.13 \pm 0.13	0.11 \pm 0.18	0.1 \pm 0.11	0.03 \pm 0.05	0.17	-
	roscoe-drop (2)	0.29 \pm 0.08	0.18 \pm 0.09	0.2 \pm 0.12	0.08 \pm 0.05	0.13 \pm 0.21	0.03 \pm 0.04	0.13	-
	roscoe-cosmos (2)	0.16 \pm 0.07	0.14 \pm 0.09	0.09 \pm 0.17	0.14 \pm 0.17	0.19 \pm 0.05	-0.03 \pm 0.01	0.13	-
	qags (1)	0.72	0.69	0.66	0.66	0.68	0.13	0.21	0.49
	medical-safety (2)	0.02 \pm 0.03	-0.02 \pm 0.02	-0.01 \pm 0.09	-0.03 \pm 0.08	0.0 \pm 0.06	0.01 \pm 0.01	0.03	-
	dices-990 (1)	-0.24	-0.16	-0.16	-0.13	-0.2	-0.09	0.05	0.14
	dices-350-expert (1)	-0.2	-0.2	-0.15	-0.03	-0.11	-0.01	0.08	-
	dices-350-crowdsourced (1)	-0.22	-0.15	-0.08	-0.02	-0.11	-0.08	0.06	0.16
	persona-chat (2)	0.24 \pm 0.34	0.15 \pm 0.21	0.58 \pm 0.59	-0.03 \pm 0.04	0.54 \pm 0.65	0.48 \pm 0.74	0.19	0.33
inferential-strategies (1)	0.42	0.08	0.02	0.22	0.06	-0.02	0.16	1.0	
Average Cohen’s κ		0.28 \pm 0.32	0.24 \pm 0.30	0.24 \pm 0.30	0.22 \pm 0.28	0.21 \pm 0.28	0.10 \pm 0.18		
Graded Annotations	dailydialog (1)	0.69	0.6	0.55	0.61	0.63	0.52	0.05	0.59
	switchboard (1)	0.66	0.6	0.63	0.59	0.56	0.36	0.13	0.57
	persona-chat (4)	0.22 \pm 0.11	0.03 \pm 0.11	0.16 \pm 0.1	0.1 \pm 0.09	0.02 \pm 0.15	0.07 \pm 0.13	0.19	0.33
	topical-chat (4)	0.26 \pm 0.03	0.21 \pm 0.15	0.13 \pm 0.04	0.17 \pm 0.12	0.21 \pm 0.18	0.14 \pm 0.05	0.15	0.08
	recipe-crowd-sourcing-data (6)	0.78 \pm 0.05	0.75 \pm 0.05	0.6 \pm 0.15	0.67 \pm 0.09	0.57 \pm 0.23	0.32 \pm 0.28	0.2	0.41
	roscoe-cosmos (2)	0.57 \pm 0.18	0.52 \pm 0.16	0.51 \pm 0.16	0.57 \pm 0.17	0.53 \pm 0.21	0.33 \pm 0.25	0.13	-
	roscoe-drop (2)	0.57 \pm 0.22	0.5 \pm 0.17	0.44 \pm 0.15	0.44 \pm 0.13	0.32 \pm 0.12	0.21 \pm 0.22	0.13	-
	roscoe-esnli (2)	0.49 \pm 0.24	0.38 \pm 0.21	0.38 \pm 0.17	0.35 \pm 0.21	0.32 \pm 0.12	0.09 \pm 0.08	0.17	-
	roscoe-gsm8k (2)	0.82 \pm 0.12	0.77 \pm 0.17	0.81 \pm 0.14	0.81 \pm 0.12	0.79 \pm 0.13	0.68 \pm 0.2	0.23	-
	newsroom (4)	0.59 \pm 0.02	0.63 \pm 0.01	0.44 \pm 0.04	0.55 \pm 0.04	0.5 \pm 0.07	0.36 \pm 0.06	0.1	0.11
	summeval (4)	0.37 \pm 0.07	0.26 \pm 0.15	0.54 \pm 0.08	0.4 \pm 0.02	0.48 \pm 0.02	0.19 \pm 0.06	0.15	-
	wmt-23-en-de (1)	0.22	0.19	0.23	0.16	0.17	0.22	0.04	-
	wmt-23-zh-en (1)	0.17	0.15	0.19	0.14	0.15	0.15	0.03	-
	wmt-human-en-de (1)	0.63	0.37	0.51	0.46	0.19	0.42	0.2	0.5
	wmt-human-zh-en (1)	0.54	0.37	0.48	0.41	0.25	0.42	0.14	0.09
	Average Spearman’s ρ		0.50 \pm 0.21	0.42 \pm 0.23	0.44 \pm 0.19	0.43 \pm 0.21	0.38 \pm 0.22	0.30 \pm 0.17	

Table 1: Scores per dataset for the models with $\geq 98\%$ valid response rates (results for all models in Table 3, App. D): Cohen’s kappa for categorical annotations and Spearman’s correlation for graded annotations. Datasets with both categorical and graded annotations appear twice. Datasets in blue concern human-generated language, while those in red concern model-generated text. ‘ σ ’ denotes the standard deviation of the scores across models per dataset (averaged over properties if more than one is judged per dataset). Krippendorff’s α in the last column.

used and show high performance across several tasks on the Open LLM and Chatbot Arena Leaderboards (Chiang et al., 2024): GPT-4o (OpenAI, 2024), LLaMA-3 (8B and 70B; AI@Meta 2024), Gemini-1.5 (Reid et al., 2024), Mixtral (8x7B and 8x22B; Jiang et al. 2024), Command R and Command R+ (Cohere and Cohere for AI, 2024a,b), OLMo (Groeneveld et al., 2024), Starling-7B (Zhu et al., 2023), and Mistral (Jiang et al., 2023a). See Appendix B for inference procedure details.

Prompts Considering that most datasets include the original instructions that were used to collect human judgments, we use them as is for the model prompts. We include additional instructions to constrain the models’ output and minimise verbosity: ‘Answer with one of {}. Do not explain your answer.’ We also experimented with system prompts but did not see any improvements. When the origi-

nal instruction to collect human judgments is not available, we derived a reasonable prompt by using relevant information from the original paper, such as the description of the task, and the definition of the metrics used for the evaluation. We provide all the prompts in the supplementary material.

Evaluation Models do not always respond to the prompts as requested (e.g., they may refuse to answer if they perceive the prompt as sensitive). We therefore use the following evaluation protocol: (i) To obtain the same number of judgments across models, we replace invalid LLM responses with random values, sampled from the set of possible classes in categorical annotations or the grade range in graded ones. Figure 4 in Appendix C shows the rate of valid responses per model. (ii) For graded annotations, we compute Spearman’s correlation between model and human judgments; for categor-

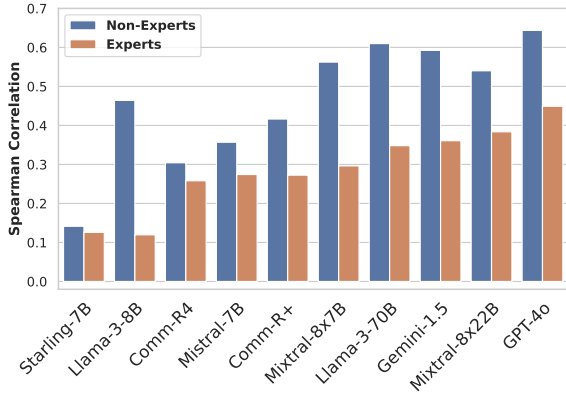


Figure 2: Model correlation with human experts vs. non-expert judgments (OLMo has a negative correlation and does not appear in the figure).

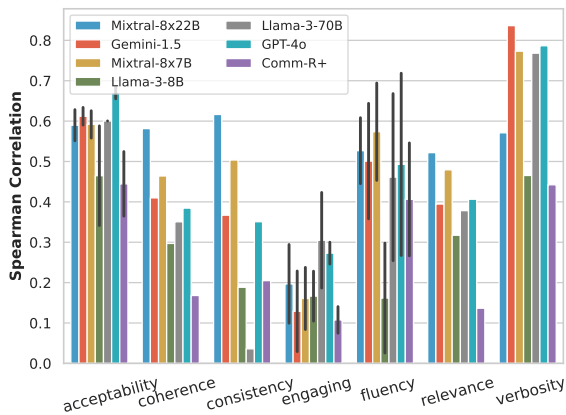


Figure 3: Correlation scores for those properties with exclusively graded judgements across datasets.

ical annotations, we compute Cohen’s κ . (iii) To get a sense for the inherent difficulty of a task, we compute human inter-rater agreement using Krippendorff’s α for the eight datasets where individual human judgments are available (Table 2, App. A).

4 Results

Our evaluation shows that scores vary substantially across models and, for any given model, they vary both across datasets and properties being judged. Table 1 presents detailed results for the 6 models that exhibit the largest rate of valid responses, i.e., $\geq 98\%$. GPT-4o ranks first across several evaluation scenarios, but the Llama-3-70B and Mistral-8x22B open models are relatively close, and outperform GPT-4o on some assessment types such as categorical sentence acceptability (CoLa) and graded summary quality (Summeval). Overall, the high degree of variability observed does not seem to be fully explained by the inherent difficulty of the annotation tasks as captured by Krippendorff’s α .

Among the property types with the lowest human-model alignment are toxicity and safety (in particular on DICES and Medical-safety), where model scores and valid response rates can be extremely low (see Figure 5 in Appendix C). One possible reason might be the RLHF guardrails associated with these tasks (Weidinger et al., 2023). We find that many models tend to provide explanations instead of outputting a judgment, especially in the medical domain.

Despite the high variability across models and datasets, we observe several interesting trends. As we can see in Figure 2 for graded annotations, all models demonstrate higher correlation results with data annotated by non-expert human judges compared to expert annotators, echoing recent findings by Aguda et al. (2024). GPT-4o exhibits an advantage over other models, which is more pronounced when evaluated against expert annotations.

Figure 3 shows correlation results across different datasets for the subset of properties that exclusively have graded judgements. We observe that the proprietary models GPT-4o and Gemini-1.5 exhibit the highest scores when evaluating acceptability and verbosity, while the two Mistral open models show the strongest correlations for coherence and consistency. Overall, no single model demonstrates a clear superiority over others across all categories; instead, different quality dimensions are better assessed by different models.

Finally, we observe that all models achieve better alignment with human judgments when evaluating human language than when assessing machine-generated text, both for categorical and graded annotations (see Figure 6 in Appendix D). This emphasises the need for caution when using LLMs to automatically evaluate the output of NLP systems.

5 Conclusions

There is an increasing trend towards turning to LLMs to automate the evaluation of model outputs, a methodology that is claimed to be a cost-efficient alternative to traditional human evaluation. We contribute to this discussion with JUDGE-BENCH, a living benchmark and a large-scale study of the correlation between human and LLM judgments across 20 datasets, considering factors such as the properties being assessed, the expertise level of the human judges, and whether the data is model- or human-generated. We find limited evidence that 11 state-of-the-art LLMs are ready to replace expert

or non-expert human judges, and caution against using LLMs for this purpose. We will make our JUDGE-BENCH available to enable future updates as new LLMs are released.

Limitations

To facilitate direct comparisons with human performance, the prompts used in our experiments were based on the original guidelines provided to the human annotators whenever these were available. However, the instruction format provided to humans may not be aligned with the format models ‘expect’ as a result of their instruction-tuning. This could limit models’ ability to provide valid and/or human-like outputs. A related limitation is that we do not take into account the valid/invalid response rates, which could affect the estimate of models’ ‘true’ correlation with human judgments. In our work, we address this limitation by replacing invalid responses with random values.

Finally, our work mostly focuses on English-language datasets – with the exception of datasets focussing specifically on machine-translation outputs. An interesting direction to explore in future work might be to check whether LLMs’ meta-evaluation abilities vary across different languages.

Acknowledgements

This work emerged from discussions at a workshop organised by RF and SP at the Oberwolfach Research Institute for Mathematics (MFO) on behalf of the ELLIS NLP programme. The event was funded by the state of Baden-Württemberg (Germany) and organised in collaboration with the ELLIS Institute Tübingen and the Max Planck Institute for Intelligent Systems. AB, EG, RF, and AT were supported by the European Research Council (ERC Consolidator Grant DREAM 819455 to RF). DE was supported by a research grant (VIL53122) from VILLUM FONDEN. MG was supported by an ETH Zurich Postdoctoral Fellowship. MH was supported in part by an OpenAI Superalignment Fellowship. AM was supported by the European Research Council (DECOLLAGE, ERC-2022-CoG 101088763) and by Fundação para a Ciência e Tecnologia through contract UIDB/50008/2020. BP was supported by an ERC Consolidator Grant (DIALECT 101043235). AKS was supported by the TIMELY project under the EU-H2020 grant 101017424.

References

- Gavin Abercrombie and Verena Rieser. 2022. [Risk-graded safety for handling medical queries in conversational AI](#). In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 234–243, Online only. Association for Computational Linguistics.
- Toyin D. Aguda, Suchetha Siddagangappa, Elena Kochkina, Simerjot Kaur, Dongsheng Wang, and Charese Smiley. 2024. [Large language models as financial data annotators: A study on effectiveness and efficiency](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 10124–10145, Torino, Italia. ELRA and ICCL.
- AI@Meta. 2024. [Llama 3 model card](#).
- Lora Aroyo, Alex Taylor, Mark Díaz, Christopher Homan, Alicia Parrish, Gregory Serapio-García, Vinodkumar Prabhakaran, and Ding Wang. 2023. [Dices dataset: Diversity in conversational ai evaluation for safety](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 53330–53342. Curran Associates, Inc.
- Simone Balloccu, Patrícia Schmidová, Mateusz Lango, and Ondrej Dusek. 2024. [Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93, St. Julian’s, Malta. Association for Computational Linguistics.
- Ipek Baris Schlicht, Defne Altiok, Maryanne Taouk, and Lucie Flek. 2024. [Pitfalls of conversational LLMs on news debiasing](#). In *Proceedings of the First Workshop on Language-driven Deliberation Technology (DELITE) @ LREC-COLING 2024*, pages 33–38, Torino, Italia. ELRA and ICCL.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31.
- Yi Chen, Rui Wang, Haiyun Jiang, Shuming Shi, and Ruifeng Xu. 2023. [Exploring the use of large language models for reference-free text quality evaluation: An empirical study](#). In *Findings of the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pages 361–374, Nusa Dua, Bali. Association for Computational Linguistics.
- Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large language models be an alternative to human evaluations?](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.

- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. 2024. Chatbot Arena: An open platform for evaluating LLMs by human preference. *arXiv preprint arXiv:2403.04132*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *ArXiv*, abs/2110.14168.
- Cohere and Cohere for AI. 2024a. [Command R Model Card](#).
- Cohere and Cohere for AI. 2024b. [Command R+ Model Card](#).
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. 2024. AlpacaFarm: A simulation framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. [SummEval: Re-evaluating Summarization Evaluation](#). *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André FT Martins, Graham Neubig, Ankush Garg, Jonathan H Clark, Markus Freitag, and Orhan Firat. 2023. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066–1083.
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. [Experts, errors, and context: A large-scale study of human evaluation for machine translation](#). *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. ChatGPT outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120.
- John J Godfrey, Edward C Holliman, and Jane McDaniel. 1992. Switchboard: Telephone speech corpus for research and development. In *Acoustics, speech, and signal processing, IEEE international conference on*, volume 1, pages 517–520. IEEE Computer Society.
- Olga Golovneva, Moya Peng Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. [ROSCOE: A suite of metrics for scoring step-by-step reasoning](#). In *The Eleventh International Conference on Learning Representations*.
- Karthik Gopalakrishnan, Behnam Hedayatnia, Qiang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tur. 2019. [Topical-chat: Towards knowledge-grounded open-domain conversations](#). In *Proc. Interspeech 2019*, pages 1891–1895.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muenhoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Nisiant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hanan Hajishirzi. 2024. Olmo: Accelerating the science of language models. *Preprint*.
- Max Grusky, Mor Naaman, and Yoav Artzi. 2018. [Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 708–719, New Orleans, Louisiana. Association for Computational Linguistics.
- Rishav Hada, Varun Gumma, Adrian Wynter, Harshita Diddee, Mohamed Ahmed, Monojit Choudhury, Kalika Bali, and Sunayana Sitaram. 2024. [Are large language model-based evaluators the solution to scaling up multilingual evaluation?](#) In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1051–1070, St. Julian’s, Malta. Association for Computational Linguistics.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28.
- Fan Huang, Haewoon Kwak, Kunwoo Park, and Jisun An. 2024. [ChatGPT rates natural language explanation quality like humans: But on which scales?](#) In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language*

- Resources and Evaluation (LREC-COLING 2024)*, pages 3111–3132, Torino, Italia. ELRA and ICCL.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. [Cosmos QA: Machine reading comprehension with contextual commonsense reasoning](#). 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 2391–2401, Hong Kong, China. Association for Computational Linguistics.
- 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. 2023a. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023b. [LLM-blender: Ensembling large language models with pairwise ranking and generative fusion](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14165–14178, Toronto, Canada. Association for Computational Linguistics.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. [Findings of the 2023 conference on machine translation \(WMT23\): LLMs are here but not quite there yet](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.
- Tom Kocmi and Christian Federmann. 2023a. [GEMBA-MQM: Detecting translation quality error spans with GPT-4](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775, Singapore. Association for Computational Linguistics.
- Tom Kocmi and Christian Federmann. 2023b. [Large language models are state-of-the-art evaluators of translation quality](#). In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.
- Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2023. [Benchmarking cognitive biases in large language models as evaluators](#). *arXiv preprint arXiv:2309.17012*.
- Fabian Landwehr, Erika Varis Doggett, and Romann M. Weber. 2023. [Memories for virtual AI characters](#). In *Proceedings of the 16th International Natural Language Generation Conference*, pages 237–252, Prague, Czechia. Association for Computational Linguistics.
- Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang. 2023. [ToxicChat: Unveiling hidden challenges of toxicity detection in real-world user-AI conversation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4694–4702, Singapore. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Shikib Mehri and Maxine Eskenazi. 2020. [USR: An unsupervised and reference free evaluation metric for dialog generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 681–707, Online. Association for Computational Linguistics.
- Philipp Mondorf and Barbara Plank. 2024. [Comparing inferential strategies of humans and large language models in deductive reasoning](#). In *The 62nd Annual Meeting of the Association for Computational Linguistics*.
- Ben Naismith, Phoebe Mulcaire, and Jill Burstein. 2023. [Automated evaluation of written discourse coherence using GPT-4](#). In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 394–403, Toronto, Canada. Association for Computational Linguistics.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. [Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.
- OpenAI. 2024. [Gpt-4o model card](#).
- Maja Pavlovic and Massimo Poesio. 2024. [The effectiveness of LLMs as annotators: A comparative overview and empirical analysis of direct representation](#). In *Proceedings of the 3rd Workshop on Perspective Approaches to NLP (NLPerspectives) @ LREC-COLING 2024*, pages 100–110, Torino, Italia. ELRA and ICCL.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.

- Ehud Reiter. 2024. [Can LLM-based eval replace human evaluation?](#) Blog post.
- Katharina Stein, Lucia Donatelli, and Alexander Koller. 2023. [From sentence to action: Splitting AMR graphs for recipe instructions](#). In *Proceedings of the Fourth International Workshop on Designing Meaning Representations*, pages 52–67, Nancy, France. Association for Computational Linguistics.
- Petter Törnberg. 2023. ChatGPT-4 outperforms experts and crowd workers in annotating political twitter messages with zero-shot learning. *arXiv preprint arXiv:2304.06588*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. [Zephyr: Direct distillation of lm alignment](#). *Preprint*, arXiv:2310.16944.
- 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*.
- Sarenne Carrol Wallbridge, Catherine Lai, and Peter Bell. 2022. [Investigating perception of spoken dialogue acceptability through surprisal](#). In *Proc. Interspeech 2022*, pages 4506–4510.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. [Asking and answering questions to evaluate the factual consistency of summaries](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5008–5020, Online. Association for Computational Linguistics.
- Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. [Is ChatGPT a good NLG evaluator? a preliminary study](#). In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11, Singapore. Association for Computational Linguistics.
- Alex Warstadt and Samuel R. Bowman. 2020. [Linguistic analysis of pretrained sentence encoders with acceptability judgments](#). *Preprint*, arXiv:1901.03438.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. [Neural Network Acceptability Judgments](#). *Transactions of the Association for Computational Linguistics*, 7:625–641.
- Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-Garcia, Stevie Bergman, Jackie Kay, Conor Griffin, Ben Bariach, et al. 2023. Sociotechnical safety evaluation of generative ai systems. *arXiv preprint arXiv:2310.11986*.
- Minghao Wu and Alham Fikri Aji. 2023. Style over substance: Evaluation biases for large language models. *arXiv preprint arXiv:2307.03025*.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2024. [Evaluating large language models at evaluating instruction following](#). In *The Twelfth International Conference on Learning Representations*.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.
- 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.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. 2023. Starling-7B: Improving LLM Helpfulness & Harmlessness with RLAIIF.

Appendix

A Datasets

This section provides brief descriptions of the datasets employed in our study. Table 2 summarises relevant dataset information. Note that dataset sizes as reported in Table 2 refer to the number of annotated samples (not to the total number of collected annotations) and might therefore differ from the figures reported in the original papers.

CoLa (Warstadt et al., 2019) The Corpus of Linguistic Acceptability (CoLA) consists of 10657 sentences from 23 linguistics publications, expertly annotated for acceptability (grammaticality) by their original authors.

CoLa grammar (Warstadt and Bowman, 2020) The dataset consists of a grammatically annotated version of the CoLA development set. Each sentence in the CoLA development set is labelled with boolean features indicating the presence or absence of a particular grammatical construction (usually

syntactic in nature). Two related sets of features are considered: 63 minor features correspond to fine-grained phenomena, and 15 major features correspond to broad classes of phenomena.

Switchboard and Dailydialog (Wallbridge et al., 2022) Switchboard includes acceptability judgements collected using stimuli from the Switchboard Telephone Corpus (Godfrey et al., 1992). More specifically, the judgements refer to how plausible it is that a specific response belongs to a telephonic dialogue. The same kind of judgements are provided for Dailydialog, which collects written dialogues intended to mimic conversations that could happen in real life.

Inferential strategies Mondorf and Plank (2024) collect annotations on the logical validity of reasoning steps that models – in this case, Llama-2-chat-hf3 (Touvron et al., 2023), Mistral-7B-Instruct-v0.2 (Jiang et al., 2023a) and Zephyr-7b-beta (Tunstall et al., 2023) – generate when prompted to solve problems of propositional logic. Binary labels are assigned to each response, indicating whether the rationale provided by the model is sound (True) or not (False). Each model is assessed on 12 problems of propositional logic across 5 random seeds, resulting in a total of 60 responses per model.

Recipe-generation (Stein et al., 2023) This dataset contains human annotations assessing the quality of machine-generated recipes based on 6 attributes: grammar, fluency, verbosity, structure, success, overall.

Medical-safety (Abercrombie and Rieser, 2022) This dataset consists of 3701 pairs of medical queries (collected from a subreddit on medical advice) and both machine-generated and human-generated answers. Queries were classified by human annotators according to their severity (from ‘Not medical’ to ‘Serious’, with ‘Serious’ indicating that emergency care would be required) and answers were categorised based on their risk level (from ‘Non-medical’ to ‘Diagnosis/Treatment’).

DICES (Aroyo et al., 2023) The DICES datasets consist of a series of machine-generated responses whose safety is judged based on the previous conversation turns (context). While the original dataset provides fine-grained annotations with answers to questions targeting specific aspects of safety, here we only consider the ‘overall’ categorisation com-

prehensive of all aspects. In DICES 990 safety is judged by crowdsourced annotators, whereas in DICES 350 both expert and crowdsourced annotations are provided.

ToxicChat (Lin et al., 2023) collect binary judgements on the toxicity and ‘jailbreaking’ nature (prompt hacks deliberately intended to bypass safety policies and induce models to generate unsafe content) of human prompts to LLMs. While the original dataset contains a mix of human- and automatically-annotated instances, here we only consider the human-annotated prompts.

Topical Chat and Persona Chat (Mehri and Eskenazi, 2020) These datasets contain human judgments on the quality of machine- and human-generated responses based on the provided dialogue context. The annotated dialogues were selected from Topical Chat (Gopalakrishnan et al., 2019) – a dataset collecting human-human conversations on provided facts – and Persona Chat (Zhang et al., 2018), which contains human-human persona-conditioned conversations. Each response is evaluated on 6 attributes: Understandable, Natural, Maintains Context, Interesting, Uses Knowledge, and Overall Quality.

WMT 2020 EnDe/ZhEn (Freitag et al., 2021) These datasets are a re-annotated version of the English-to-German and Chinese-to-English test sets taken from the WMT 2020 news translation task. The annotation was carried out by raters who are professional translators and native speakers of the target language using a Scalar Quality Metric (SQM) evaluation on a 0–6 rating scale.

WMT 2023 EnDe/ZhEn (Kocmi et al., 2023) These datasets are the English-to-German and Chinese-to-English test sets taken from the General Machine Translation Task organised as part of the 2023 Conference on Machine Translation (WMT). In contrast to previous editions, the evaluation of translation quality was conducted by a professional or semi-professional annotator pool rather than utilising annotations from MTurk. Annotators were asked to provide a score between 0 and 100 on a sliding scale.

ROSCOE (Golovneva et al., 2023) collect human judgments assessing the quality of GPT-3’s reasonings. The output reasonings are elicited by inputting GPT-3 with questions selected from 4 commonly used reasoning datasets, i.e.,

Dataset	Task	Size	Type	Guidelines	Expert	Agreement	Leaked
CoLA (Warstadt et al., 2019)	Acceptability	1,043	Categorical	✗	✓	✓	✓
CoLA grammar (Warstadt and Bowman, 2020)	Acceptability	1,043	Categorical	✗	✓	✗	✓
Switchboard (Wallbridge et al., 2022)	Acceptability	100	Graded	✓	✗	✗	
Dailydialog (Wallbridge et al., 2022)	Acceptability	100	Graded	✓	✗	✗	
Inferential strategies (Mondorf and Plank, 2024)	Reasoning	300	Categorical	✓	✓	✗	✗
ROSCOE (Golovneva et al., 2023)	Reasoning	756	Categorical + Graded	✓	✓	✗	
Recipe-generation (Stein et al., 2023)	Planning	52	Graded	✓		✗	
Medical-safety (Abercrombie and Rieser, 2022)	Toxicity & Safety	3,701	Preference	✓	✓		
DICES (Aroyo et al., 2023)	Toxicity & Safety	1,340	Categorical	✗	Mixed	✓	
ToxicChat (Lin et al., 2023)	Toxicity & Safety	5,654	Categorical	✗	✓	✗	
Topical Chat (Mehri and Eskenazi, 2020)	Dialogue	60	Graded + Categorical	✗	✓	✓	
Persona Chat (Mehri and Eskenazi, 2020)	Dialogue	60	Graded + Categorical	✗	✓	✓	
WMT 20 EnDe (Freitag et al., 2021)	Machine Translation	14,122	Graded	✗	✓	✓	
WMT 20 ZhEn (Freitag et al., 2021)	Machine Translation	19,974	Graded	✗	✓	✓	
WMT 23 EnDe (Kocmi et al., 2023)	Machine Translation	6,588	Graded	✗	✓	✗	
WMT 23 ZhEn (Kocmi et al., 2023)	Machine Translation	13,245	Graded	✗	✓	✗	
G-Eval / SummEval (Liu et al., 2023)	Summarisation	1,600	Graded	✓		✓	✓
QAGS (Wang et al., 2020)	Summarisation	953	Categorical	✓	✗	✗	
NewsRoom (Grusky et al., 2018)	Summarisation	420	Graded	✓	✗	✗	✓
LLMBar (Zeng et al., 2024)	Instruction Following	419	Categorical	✓	✓	✓	✗

Table 2: Overview of the main features of the datasets considered in the study. Note that ‘Size’ refers to the number of annotated samples, not to the total number of human annotations. ‘Agreement’ indicates whether multiple annotations were available for the same instance or not. Information on possible data leakage was retrieved from Balloccu et al. (2024).

CosmosQA (Huang et al., 2019), DROP (Dua et al., 2019), e-SNLI (Camburu et al., 2018) and GSM8K (Cobbe et al., 2021). While ROSCOE provides annotations on each step of the reasoning trace, here we only consider the global judgments over the whole reasoning.

G-Eval/Summeval (Liu et al., 2023; Fabbri et al., 2021) These datasets include summaries generated by multiple recent summarisation models trained on the CNN/DailyMail dataset (Hermann et al., 2015). Summaries are annotated by both expert judges and crowdsourced workers on 4 dimensions: coherence, consistency, fluency, relevance.

QAGS (Wang et al., 2020) QAGS consists of annotations judging the factual consistency of one-sentence model-generated summaries of news articles. The gold-standard summaries and articles are collected from CNN/DailyMail (Hermann et al., 2015) and XSUM (Narayan et al., 2018).

NewsRoom (Grusky et al., 2018) This dataset includes human judgments on the quality of system-generated summaries of news articles. More specifically, annotators evaluated summaries across two semantic dimensions (informativeness and relevancy) and two syntactic dimensions (fluency and coherence).

LLMBar (Zeng et al., 2024) LLMBar is a dataset targeted at evaluating the instruction-following abilities of LLMs. Each entry of this dataset consists of an instruction paired with two

different outputs, one correctly following the instruction and the other deviating from it. LLMBar has an adversarial split where deviating outputs are carefully constructed to ‘fool’ LLM-based evaluators and a natural split where deviating outputs are more naturalistic.

B Inference Details

All open-model checkpoints were obtained using the HuggingFace pipeline and we access all proprietary models using their corresponding API libraries. The proprietary models were accessed between 06-06-2024 and 13-06-2024. We obtain the model responses using greedy decoding, which we operationalize for the proprietary models by setting the temperature parameter to 0. We allow open models to generate a maximum of 25 new tokens, and proprietary models to generate a maximum of 5 new tokens.

We leverage Nvidia A100 (80 GB) GPUs for a total of 125.22 compute hours. The cost of running experiments on gemini-1.5-flash was €20.65, while the cost of experiments on GPT-4o was approximately \$100.

C Models’ Valid Response Rates

Figures 4 and 5 show the rate of valid responses per model and per dataset, respectively.

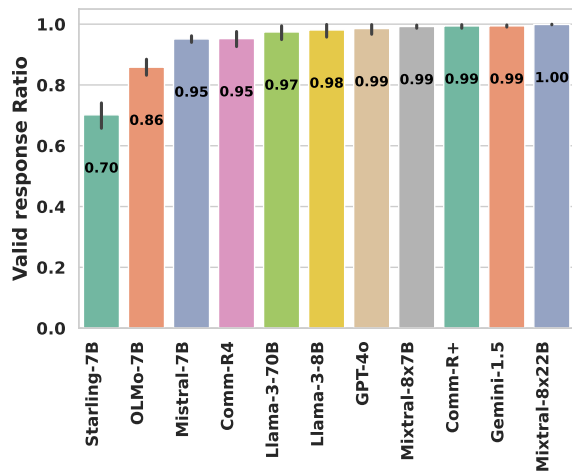


Figure 4: Valid response rate per model.

D Additional Results

In Table 3 we report human-model alignment scores per dataset for all models tested, thus complementing Table 1 in the paper. Figure 6 shows alignment scores broken down according to the source of the material to be judged, i.e., human or machine generated output.

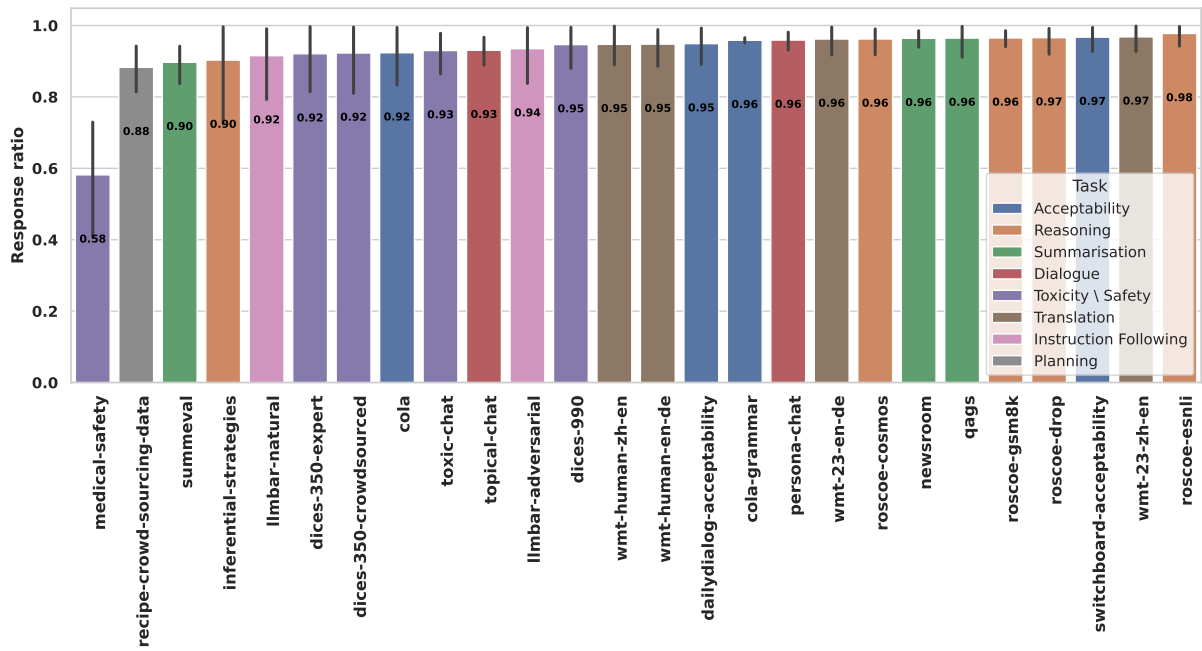


Figure 5: Average ratios of valid responses across datasets over the 11 models we tested.

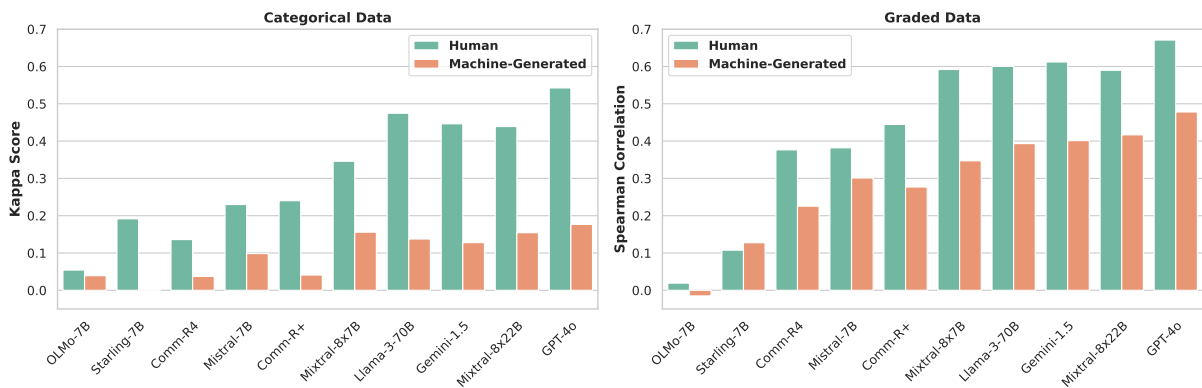


Figure 6: Scores (Cohen’s κ for categorical annotations and Spearman’s correlation for graded annotations) on test items involving human language vs. machine-generated outputs.

Type	Dataset (# properties judged)	GPT-4o	Llama-3-70B	Mixtral-8x22B	Gemini-1.5	Mixtral-8x7B	Comm-R+	Comm-R4	Llama-3-8B	Mistral-7B	Starling-7B	OLMo-7B	
Categorical Annotations	cola (1)	0.34	0.56	0.54	0.44	0.55	0.12	0.01	0.49	0.35	0.45	0.05	
	limbar-natural (1)	0.84	0.72	0.72	0.79	0.5	0.56	0.62	0.52	0.38	0.03	-0.03	
	limbar-adversarial (1)	0.58	0.4	0.2	0.29	0.07	0.11	-0.19	-0.17	-0.07	-0.19	0.08	
	toxic-chat (2)	0.49 ±0.36	0.4 ±0.13	0.45 ±0.27	0.45 ±0.35	0.36 ±0.12	0.28 ±0.35	0.2 ±0.21	0.4 ±0.27	0.4 ±0.14	0.06 ±0.08	0.27 ±0.26	-0.0 ±0.03
	cola-grammar (63)	0.47 ±0.22	0.29 ±0.24	0.28 ±0.23	0.26 ±0.24	0.21 ±0.18	0.13 ±0.14	0.08 ±0.1	0.17 ±0.18	0.06 ±0.08	0.07 ±0.08	0.01 ±0.02	0.01 ±0.02
	topical-chat (2)	0.05 ±0.07	-0.02 ±0.02	-0.03 ±0.04	-0.03 ±0.04	0.02 ±0.03	0.01 ±0.02	0.01 ±0.01	0.69 ±0.44	0.06 ±0.08	0.04 ±0.06	0.04 ±0.06	0.04 ±0.05
	roscoe-gsm8k (2)	0.59 ±0.35	0.65 ±0.27	0.62 ±0.38	0.6 ±0.24	0.58 ±0.36	0.0	0.21 ±0.03	0.18 ±0.24	0.56 ±0.42	-0.02 ±0.01	-0.02 ±0.01	-0.02 ±0.08
	roscoe-esnli (2)	0.29 ±0.06	0.31 ±0.15	0.13 ±0.13	0.11 ±0.18	0.1 ±0.11	0.03 ±0.05	-0.01 ±0.01	0.02 ±0.04	0.11 ±0.02	0.01 ±0.03	0.03 ±0.04	-0.02 ±0.03
	roscoe-drop (2)	0.29 ±0.08	0.18 ±0.09	0.2 ±0.12	0.08 ±0.05	0.13 ±0.21	0.03 ±0.04	0.02 ±0.07	0.01 ±0.03	0.08 ±0.01	0.03 ±0.04	0.03 ±0.04	-0.07 ±0.03
	roscoe-cosmos (2)	0.16 ±0.07	0.14 ±0.09	0.09 ±0.17	0.14 ±0.17	0.19 ±0.05	-0.03 ±0.01	-0.01 ±0.02	0.07 ±0.09	0.22 ±0.06	-0.01 ±0.01	-0.01 ±0.01	-0.07 ±0.03
	gags (1)	0.72	0.69	0.66	0.66	0.68	0.13	0.33	0.52	0.22	0.02	0.02	0.35
	medical-safety (2)	0.02 ±0.03	-0.02 ±0.02	-0.01 ±0.09	-0.03 ±0.08	0.0 ±0.06	0.01 ±0.01	0.02 ±0.01	-0.01 ±0.02	-0.02 ±0.07	0.01 ±0.01	0.01 ±0.01	0.04 ±0.02
	medical-990 (1)	-0.24	-0.16	-0.16	-0.13	-0.2	-0.09	-0.02	-0.12	-0.01	-0.01	-0.05	-0.02
	dices-350-expert (1)	-0.2	-0.2	-0.15	-0.03	-0.11	-0.01	0.01	-0.06	0.01	0.01	0.0	0.03
	dices-350-crowdsourced (1)	-0.22	-0.15	-0.08	-0.02	-0.11	-0.08	0.01	-0.07	-0.01	-0.04	-0.04	0.04
	persona-chat (2)	0.24 ±0.34	0.15 ±0.21	0.58 ±0.59	-0.03 ±0.04	0.54 ±0.65	0.48 ±0.74	0.01 ±0.01	0.47 ±0.74	0.0 ±0.01	0.09 ±0.13	0.04 ±0.05	0.04 ±0.05
	inferential-strategies (1)	0.42	0.08	0.02	0.22	0.06	-0.02	-0.12	-0.01	-0.03	-0.01	-0.01	0.03
	Graded Annotations	dailydiarlog (1)	0.69	0.6	0.55	0.61	0.63	0.52	0.26	0.59	0.36	0.09	0.02
		switchboard (1)	0.66	0.6	0.63	0.59	0.56	0.36	0.58	0.34	0.4	0.01	-0.04
		persona-chat (4)	0.22 ±0.11	0.03 ±0.11	0.16 ±0.1	0.1 ±0.09	0.02 ±0.15	0.07 ±0.13	0.06 ±0.2	0.01 ±0.15	-0.08 ±0.12	0.03 ±0.18	0.03 ±0.16
topical-chat (4)		0.26 ±0.03	0.21 ±0.15	0.13 ±0.04	0.17 ±0.12	0.21 ±0.18	0.14 ±0.05	0.06 ±0.07	0.19 ±0.11	0.25 ±0.21	0.13 ±0.1	0.01 ±0.12	
recipe-crowd-sourcing-data (6)		0.78 ±0.05	0.75 ±0.05	0.6 ±0.15	0.67 ±0.09	0.57 ±0.23	0.32 ±0.28	0.07 ±0.2	0.38 ±0.17	0.5 ±0.14	0.05 ±0.19	-0.13 ±0.13	
roscoe-cosmos (2)		0.57 ±0.18	0.52 ±0.16	0.51 ±0.16	0.57 ±0.17	0.53 ±0.21	0.33 ±0.25	0.49 ±0.16	0.14 ±0.17	0.56 ±0.24	0.11 ±0.02	0.13 ±0.08	
roscoe-drop (2)		0.57 ±0.22	0.5 ±0.17	0.44 ±0.15	0.44 ±0.13	0.32 ±0.12	0.21 ±0.22	0.37 ±0.18	0.22 ±0.11	0.36 ±0.16	0.18 ±0.15	0.07 ±0.1	
roscoe-esnli (2)		0.49 ±0.24	0.38 ±0.21	0.38 ±0.17	0.35 ±0.21	0.32 ±0.12	0.09 ±0.08	0.28 ±0.21	-0.07 ±0.15	0.23 ±0.23	0.13 ±0.09	-0.01 ±0.12	
roscoe-gsm8k (2)		0.82 ±0.12	0.77 ±0.17	0.81 ±0.14	0.81 ±0.12	0.79 ±0.13	0.68 ±0.2	0.7 ±0.07	0.34 ±0.06	0.43 ±0.11	0.5 ±0.12	-0.15 ±0.07	
newstroom (4)		0.59 ±0.02	0.63 ±0.01	0.44 ±0.04	0.55 ±0.04	0.5 ±0.07	0.36 ±0.06	0.16 ±0.05	0.46 ±0.12	0.31 ±0.05	0.18 ±0.08	-0.02 ±0.05	
summeval (4)		0.37 ±0.07	0.26 ±0.15	0.54 ±0.08	0.4 ±0.02	0.48 ±0.02	0.19 ±0.06	0.12 ±0.07	0.21 ±0.12	0.37 ±0.06	0.15 ±0.06	0.0 ±0.03	
wmt-23-en-de (1)		0.22	0.19	0.23	0.16	0.17	0.22	0.19	0.13	0.16	-0.09	0.03	
wmt-23-zh-en (1)		0.17	0.15	0.19	0.14	0.15	0.15	0.13	0.08	0.14	0.01	-0.03	
wmt-human-en-de (1)		0.63	0.37	0.51	0.46	0.19	0.42	0.15	0.05	0.31	0.15	-0.05	
wmt-human-zh-en (1)		0.54	0.37	0.48	0.41	0.25	0.42	0.15	0.12	0.38	0.14	-0.01	

Table 3: Scores per dataset for all models we evaluate. Cohen’s kappa for categorical annotations and Spearman’s correlation for graded annotations. Datasets in blue concern human-generated language while those in red concern model-generated text.