Cross-Lingual Auto Evaluation for Assessing Multilingual LLMs

Sumanth Doddapaneni*^{1,2} Mohammed Safi Ur Rahman Khan*^{1,2} Dilip Venkatesh¹ Raj Dabre^{1,2,4} Anoop Kunchukuttan^{1,3} Mitesh M. Khapra^{1,2}

¹Nilekani Centre at AI4Bharat ²IIT Madras ³Microsoft ⁴National Institute of Information and Communications Technology, Kyoto, Japan

Correspondence: {sumanthd, miteshk}@cse.iitm.ac.in, safikhan@ai4bharat.org

huggingface.co/CIA-Suite

github.com/CIA

Abstract

Evaluating machine-generated text remains a significant challenge in NLP, especially for non-English languages. Current methodologies, including automated metrics, human assessments, and LLM-based evaluations, predominantly focus on English, revealing a significant gap in multilingual evaluation frameworks. We introduce the Cross Lingual Auto Evaluation (CIA) Suite, an extensible framework that includes evaluator LLMs (HERCULE) and a novel test set (RECON) specifically designed for multilingual evaluation. Our test set features 500 human-annotated instructions spanning various task capabilities along with human judgment scores across six languages. This would enable benchmarking of general-purpose multilingual LLMs and facilitate meta-evaluation of Evaluator LLMs. The proposed model, HERCULE, is a cross-lingual evaluation model that addresses the scarcity of reference answers in the target language by learning to assign scores to responses based on easily available reference answers in English. Our experiments demonstrate that HERCULE aligns more closely with human judgments compared to proprietary models, demonstrating the effectiveness of such cross-lingual evaluation in low resource scenarios. Further, it is also effective in zero-shot evaluation on unseen languages. This study is the first comprehensive examination of crosslingual evaluation using LLMs, presenting a scalable and effective approach for multilingual assessment. All code, datasets, and models will be publicly available to enable further research in this important area.

1 Introduction

Evaluating machine-generated text has long been a central challenge in Natural Language Processing (NLP). Substantial progress has been made in English-language evaluations, using (i) automated metrics (Papineni et al., 2002; Lin, 2004;

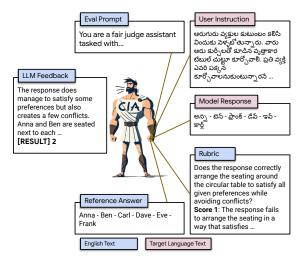


Figure 1: We present cross-lingual Evaluator LLM, HERCULE, where the **Instruction** & **Response** provided to the model are in the target language, while all other fields are in English. The model generates feedback & score in English for a given evaluation example.

Rei et al., 2020; Sellam et al., 2020), (ii) human evaluations (Watts et al., 2024; Chiang et al., 2024), and (iii) more recently, automated evaluations using Large Language Models (LLMs) (Zheng et al., 2023a; Kim et al., 2023b; Dubois et al., 2023; Li et al., 2023). However, a comprehensive framework for multilingual evaluation which goes beyond English remains missing. This gap is largely due to the absence of a robust multilingual benchmark covering complex open-ended tasks and a reliable evaluation metric. Time and again, benchmarks have proven essential in driving progress in NLP (Rajpurkar et al., 2016; Hu et al., 2020a; Wang et al., 2018; Conneau et al., 2018). This leads us to our first observation: $\{A\}$ there is an urgent need to develop a robust multilingual benchmark and an easy-to-use evaluation framework to further advance the field.

Even if one develops a comprehensive benchmark that encompasses a wide variety of gener-

^{*}Equal Contribution.

ation tasks and instruction-following capabilities, the challenge of achieving quick and automated evaluation persists. Traditionally, human evaluations have been the most reliable method for assessing models. However, as these models become more sophisticated and creative, it has become difficult for non-experts to evaluate model outputs accurately. Many non-experts often depend on superficial indicators of correctness, resulting in human evaluations that have devolved into mere "vibe checks" (Chiang et al., 2024), wherein, evaluators rely on personal biases rather than objective criteria to determine a winner. This limitation has led to the growing adoption of LLMs as evaluators, given that LLMs often possess more extensive knowledge bases than human evaluators. While LLM-based evaluation comes with its own challenges (Doddapaneni et al., 2024; Zeng et al., 2023; Zheng et al., 2023a), LLMs offer faster, cheaper, and, in certain contexts, more reliable assessments compared to human evaluations.

Grounded in this modern reality, we focus on a framework for LLM based cross-lingual evaluation. To begin with, we make a few observations. First, prior works (Kim et al., 2023b, 2024b; Wang et al., 2023b) have shown that trained evaluators significantly outperform untrained ones, matching the performance of proprietary LLMs. $\{\mathcal{B}\}\$ Hence, it is prudent to build a trained crosslingual LLM evaluator. Second, within the space of trained evaluators, some studies (Doddapaneni et al., 2024; Kim et al., 2023a) have shown that referencebased approaches (Kim et al., 2023a) where the evaluator model is provided with a reference answer are more accurate and reliable than referencefree approaches (Hada et al., 2023; Zheng et al., 2023b). The latter solely rely on their parameterized knowledge and hence do not produce very reliable evaluations. $\{C\}$ We, thus focus on building a reference-based evaluator LLM. However, such references are scarcely available for non-English languages but available in abundance for English. $\{\mathcal{D}\}\$ We thus make a case for cross-lingual evaluation, wherein responses generated in a non-English language are assessed using a reference answer available in English.

Based on \mathcal{A} , \mathcal{B} , \mathcal{C} and \mathcal{D} , we introduce the Cross LIngual Auto Evaluation (CIA) Suite, an extensible framework of evaluator LLMs and datasets designed for multilingual evaluation. In this setup,

both the questions and responses are provided in the target language, while the reference answers, evaluation instructions, and rubrics remain in English, facilitating cross-lingual evaluation (ref. Fig. 1). We also introduce the **RECON** test set, ahuman-annotated, multi-purpose benchmark spanning six languages: Bengali, German, French, Hindi, Telugu, & Urdu. This test set aims to benchmark general-purpose multilingual LLMs across various tasks and to meta-evaluate the Evaluator LLMs. We construct the INTEL training set by automatically translating the Feedback-Collection dataset (Kim et al., 2023b). Finally, we release **HERCULE**, A a series of cross-lingual, referencebased Evaluator LLMs fine-tuned on INTEL using the Llama-3.1-8B model series.

We present the results of our evaluation on the RECON test set, highlighting the improved performance of fine-tuned models over their zero-shot counterparts (Sec. §5.1). Our findings demonstrate that models trained on INTEL not only outperform large, proprietary LLMs but also exhibit greater alignment with human judgments, particularly in low-resource languages (§5.2). The ablation studies show that models trained on one language can effectively perform zero-shot evaluations on others (§6.1), demonstrating the potential of cross-lingual transfer. We also assess the impact of reference answers (§6.2). We show that training from an instruction-tuned model accelerates convergence. We also highlight the benefits of LoRA training (§6.3). Finally, we use weight merging techniques to create a unified Evaluator LLM for all target languages (§6.4). All artifacts will be made publicly available.

2 Related Work

LLMs as Evaluators. With the growing openended generation capabilities of LLMs (Meta, 2024; Team et al., 2024a; Zhao et al., 2024) and the challenges of costly, inconsistent human evaluations (Chiang and Lee, 2023; Chen et al., 2023), many studies now propose using LLMs to score model outputs (Zheng et al., 2023a; Dubois et al., 2023; Kim et al., 2023b). This approach is gaining traction due to the absence of reliable task-specific metrics and strong correlations between LLM and human scores (Dubois et al., 2023; Zheng et al., 2023a). Evaluation strategies with LLMs fall into two main types: absolute grading (Liu et al., 2023; Hada et al., 2023) and pairwise com-

¹hf.co/blog/clefourrier/llm-evaluation

parison (Wang et al., 2023a; Zheng et al., 2023a), with methods relying on either prompt-based evaluations (Zheng et al., 2023a) or fine-tuning models for evaluation (Kim et al., 2023b; Wang et al., 2023b). These methods operate in both referencedriven (Fu et al., 2023; Kim et al., 2023b) and reference-free scenarios (Liu et al., 2023). Studies show that reference-based approaches are generally more reliable (Doddapaneni et al., 2024; Kim et al., 2023b), and trained evaluators demonstrate better task adaptability and correlation with human judgments compared to those relying on parametric knowledge (Kim et al., 2024b). Advanced techniques also explore multi-agent interactions (Chan et al., 2023) or external evaluators (Min et al., 2023). While LLM-based evaluation has its limitations (Doddapaneni et al., 2024; Zeng et al., 2023), it remains the preferred approach due to the lack of scalable, cost-effective alternatives.

Multilingual Evaluators. The lack of comprehensive benchmarks and reliable evaluation methods hinders progress in multilingual model development (Hu et al., 2020b; Doddapaneni et al., 2023). Existing multilingual benchmarks are limited in scale, domain, and rely heavily on costly human evaluations, focusing mainly on classification and sentence generation tasks (Singh et al., 2024; Watts et al., 2024; Doddapaneni et al., 2023; Ahuja et al., 2024). Though human evaluations and Elo ratings help create leaderboards, a reliable metric for iterative model development remain missing. Additionally, prior work shows that GPT-4, as a multilingual evaluator, delivers inconsistent results, highlighting the need for a robust multilingual benchmark and evaluation framework (Hada et al., 2023).

Weight Merging. Weight merging has proven effective in creating unified models and improving performance across tasks such as language modeling (Matena and Raffel, 2022; Li et al., 2022; Ilharco et al., 2023), instruction following (Jang et al., 2023b; Yu et al., 2024), preference learning (Jang et al., 2023a; Ramé et al., 2023), and multilingual applications (Chronopoulou et al., 2023). Techniques like linear merging (Wortsman et al., 2022) average model weights, while task vector arithmetic (Ilharco et al., 2023) uses element-wise subtraction to represent fine-tuned models and TIES (Yadav et al., 2023) resolves interference. These methods are increasingly popular for building unified multitask models. In this work, we explore merging techniques to develop a unified model capable of evaluating multiple languages.

3 CIA: Cross Lingual Auto Evaluation

We introduce the CIA suite, a comprehensive framework for cross-lingual evaluation. In this setup, the questions and responses are provided in the target language, while the reference answers, evaluation instructions, and scoring rubrics remain in English, enabling effective cross-lingual evaluation (ref. Fig. 1). This section presents the human-annotated RECON test set (§3.1), details the training data used in our experiments (§3.2), and describes the HERCULE evaluator models trained on this data (§3.3).

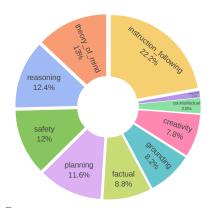


Figure 2: Distribution of task capabilities in RECON.

3.1 RECON: Test Data

We introduce RECON, a human-annotated, generalpurpose multilingual evaluation benchmark. The input prompts are fully human-generated with multiple levels of supervision. This benchmark serves two key purposes: (i) to assess the multilingual capabilities of LLMs and (ii) to meta-evaluate the performance of Evaluator LLMs. Each instance in RECON consists of a tuple $(P^X, C^{En}, R^X_{eval}, R^{En}_{ref})$ s), with superscripts indicating the language. Here, P^{X} represents the input prompt, R_{ref}^{En} denotes the reference response, C^{En} defines the evaluation criteria and rubrics. Further, for meta-evaluation of Evaluator LLMs, ${\cal R}^{X}_{eval}$ is provided as the response with an associated expected score (s). To ensure the integrity of the benchmark, all components undergo thorough manual review and validation.

The input prompts (P) are sourced from various benchmarks, ensuring all are human-written and reflective of real-world scenarios (see Fig. 2). The benchmark consists of 500 prompts, of which 250 are from BigGenBench (Kim et al., 2024a), which includes per-example rubrics and tasks like planning, instruction following, and reasoning. The

remaining 250 are curated from test sets such as UltraEval (Ding et al., 2023), WizardLM (Xu et al., 2023), LIMA (Zhou et al., 2023), MT Bench (Zheng et al., 2023a), and FBI (Doddapaneni et al., 2024), covering tasks like long-form writing, creativity, and factual questions.

Following an approach similar to BigGen-Bench (Kim et al., 2024a) and Prometheus (Kim et al., 2023b), we generate scoring criteria C for each question. We prompt GPT-40 to first generate question-specific criteria that could be used to evaluate a response, followed by detailed rubrics for assigning scores from 1 to 5. To guide the process, we provide three manually written criteria and rubric examples as in-context demonstrations. All generated rubrics are manually reviewed and verified by the authors.

To generate reference answers R_{ref} , we prompt GPT-40 with the question and corresponding rubric, instructing it to produce an answer that scores 5, formally expressed as $f(P,C,s=5) \rightarrow R_{ref}$. A manual review of sampled references confirmed their high accuracy. Similarly, we generate evaluation responses R_{eval} by prompting GPT-40 to produce answers scoring 'i' based on the rubrics, ensuring a uniform score distribution across the benchmark. This process is represented as $f(P,C,R_{ref},s=i) \rightarrow R_{eval}$, where $1 \le i \le 5$. All responses are manually verified by the authors to ensure they align with the intended score. Detailed prompts are provided in Appendix C.

Finally, all the prompts (P) in the test set are manually translated into six target languages: Bengali (bn), German (de), French (fr), Hindi (hi), Telugu (te), and Urdu (ur), with one dedicated inhouse translator assigned to each language. Additionally, the responses to be evaluated (R_{eval}) are translated using GPT-40, followed by thorough human verification and correction. Annotators were specifically instructed to ensure that errors in low-scored answers were accurately translated without any unintended corrections.

3.2 INTEL: Training Data

We use the Feedback-Collection (Kim et al., 2023b) dataset for training and the Feedback-Bench (Kim et al., 2023b) dataset for validation. During training and validation, prompts (P) and answers (R_{eval}) are translated into the target languages, while evaluation instructions, rubrics (C), and reference an-

swers (R_{ref}) remain in English. Using GPT-40,² we translate the data into six target languages, relying on automated translation due to the impracticality of manual translation.

To assess translation quality, we sampled 100 examples per language and had human experts evaluate them on a binary scale. Translations were marked valid if they conveyed the intended meaning without major errors. We found that fewer than 5% samples were deemed invalid (See Appendix D). Based on this feedback, we decided to rely on GPT-4o for our use case. Notably, this auto-translated data is used only for training and is not part of the manually verified RECON benchmark. The final dataset after auto-translation includes 100k training samples and 1k validation samples per language, formatted as $\{P^X, C^{En}, R_{ref}^{En}, R_{eval}^X, F^{En}, s\}$, where (F, s) represents the feedback and score.

3.3 HERCULE: Fine-tuned Evaluator LLM

Using INTEL (§3.2), we fine-tune LLAMA-3.1-8B-I (Meta, 2024) to equip the model with evaluation capabilities, resulting in the creation of the HERCULE series for all six target languages. We train HERCULE on the absolute grading objective, where, given a prompt (P^X) , an LLM response (R_{eval}^X) in the target language, an evaluation criteria (C^{En}) , and a reference answer in English (R_{ref}^{En}) , the evaluator LLM is tasked with providing feedback (F^{En}) and assigning a score (s) on a Likert scale from 1 to 5. Formally, we denote this objective as $f(P^X, C^{En}, R_{ref}^{En}, R_{eval}^X) \rightarrow (F^{En}, s)$. Building on prior work (Wei et al., 2022; Kim et al., 2023a), we train the model to first generate an explanation for the evaluation, followed by a score.

4 Experimental Setup

Our goal is to use our RECON test set to assess the utility of HERCULE and other LLM-based evaluators for crosslingual evaluation. In this section, we outline our experimental setup for doing so. We begin by outlining the training details (§4.1), metrics (§4.2), followed by a description of the various models considered in our experiments (§4.3).

4.1 Training Details

The HERCULE model is trained on INTEL with a sequence length of 4096, using FlashAttention 2 (Dao, 2024) and optimized with AdamW at a

²https://openai.com/index/hello-gpt-4o/

1e-5 learning rate for 3 epochs. All experiments were run on 8 Nvidia H100 GPUs. Evaluations are performed on RECON (§3.1), using human-written prompts reflecting real-world scenarios (see Fig. 2).

4.2 Metrics

As mentioned earlier, every response in the RECON test set has a ground truth score associated with it. To assess the agreement between the ground truth scores and the scores assigned by the evaluator LLM, we use linear weighted Cohen's Kappa (κ) (Sakai, 2021). An agreement score approaching 1 indicates a strong correlation between the evaluator and the ground truth, while scores approaching 0 indicate a weak alignment.

4.3 Models Considered

We consider the following models for both Zero-Shot and Fine-Tuning experiments. For fine-tuning, we train GEMMA, SARVAM, and AYA23 using the same setup as HERCULE. In the zero-shot setting, we evaluate LLAMA-3.1-405B-I, GPT-40, and GEMINI-1.5-PRO, as their large size and proprietary nature limit fine-tuning.

- **G**EMMA: An open-source LLM trained on 6T tokens with better tokenizer fertility compared to LLAMA-3.1 (Team et al., 2024b).
- SARVAM³: An open-source LLM trained for Indian languages with 2T synthetic tokens.
- ≺AYA23: Aya23 is an 8B open-weight instruction fine-tuned model, showcasing highly advanced multilingual capabilities in 23 languages (Aryabumi et al., 2024).
- [∞]LLAMA-3.1-405B-I: One of the largest opensource models currently available, evaluated in a zero-shot setting on RECON (§3.1) benchmark (Dubey et al., 2024).
- ◆GEMINI-1.5-PRO: A powerful proprietary LLM, known for its advanced multilingual capabilities compared to other closed-source models (Team et al., 2024a).

5 Results

In this section, we analyze the results of the CIA framework on the RECON test set (§5.1). We then compare its evaluation capabilities against human assessments (§5.2) and conclude with a qualitative evaluation of the LLM evaluations (§5.3).

5.1 Does Cross Lingual Evaluation Work?

We evaluate all models and baselines on the RE-CON test set, reporting the Cohen's Kappa (κ) score (Sakai, 2021) between the Evaluator LLM and the ground truth scores in Table 1. Our results demonstrate that models fine-tuned with IN-TEL consistently outperform their zero-shot counterparts. Notably, models trained on INTEL significantly surpass even large, proprietary, black-box LLMs in performance (e.g., our HERCULE model outperforms the GPT-40 model). It is important to note that, despite some models having high fertility for the languages being trained and evaluated (ref. Fig. 4 in Appendix B), they still significantly outperform zero-shot evaluations conducted with large LLMs. This emphasizes that, even when the base model does not have a fair representation for the languages of interest, fine-tuned models demonstrate superior alignment and performance compared to generic large-scale models that haven't been specifically trained on the evaluation task.

5.2 Evaluation in the Wild

So far, we evaluated the models using responses from RECON (Sec. §3.1), constructed in a controlled setting. To further assess the Evaluator LLMs in real-world scenarios, we conducted human evaluations. We sampled 100 prompts per language from RECON and generated responses using LLAMA-3.1-8B-I, GEMMA-2-2B, and GPT-40-MINI. Native speakers with formal training in their respective languages scored each response, with the final human score averaged across three annotators. We then evaluate these responses using GPT-40, GEMINI-1.5-PRO, LLAMA-3.1-405B-I, and HERCULE and compared their correlations with human judgments. On average, HERCULE demonstrated stronger alignment with human evaluations than both GPT-40 and GEMINI-1.5-PRO. For high-resource languages like hi, which are officially supported by these models (OpenAI, 2024; Meta, 2024; Aryabumi et al., 2024), zero-shot GPT-40 achieved the highest performance. However, for other languages, our cross-lingual fine-tuning provided clear advantages. We also found reasonable inter-annotator agreement (IAA in Table 2). Note that, even in this in-the-wild setup, the Evaluator LLMs maintain the same relative rankings as in RE-CON, validating RECON as a reliable benchmark for evaluating Evaluator LLMs. We report other metrics, including Kendall's Tau (τ) and Spearman's

³huggingface.co/models/sarvamai/sarvam-2b-v0.5

Model	Туре	bn	de	fr	hi	te	ur	avg.
	Zero-Shot	0.64	0.66	0.65	0.64	0.61	0.64	0.64
◆ Gemini-1.5-Pro	Zero-Shot	0.54	0.58	0.59	0.57	0.53	0.57	0.56
∞ Llama-3.1-405B-I	Zero-Shot	0.60	0.66	0.66	0.62	0.51	0.65	0.62
∞LLAMA-3.2 3B	FFT	0.68	0.72	0.71	0.71	0.70	0.72	0.71
G Gemma 7B	FFT	0.47	0.39	0.36	0.43	0.33	0.38	0.39
≺ AYA23 8B	FFT	0.70	0.72	0.73	0.72	0.65	0.71	0.70
Mercule 8B	FFT	0.74	0.75	0.75	0.74	0.69	0.74	0.73
↑ HERCULE 8B	LoRA	0.72	0.74	0.72	0.72	0.70	0.70	0.72

Table 1: Evaluation results of all models on the RECON test set. We report the Linear Weighted Cohen's Kappa (κ) score between the ground truth scores and the model predictions. Higher the value, better is the correlation. The upper half of the table presents zero-shot evaluations, while the lower half shows the results of fine-tuned models. Refer to Sec. §5.1 for detailed results.

Model	bn	hi	te	ur
	0.37	0.61	0.62	0.67
◆ GEMINI-PRO	0.31	0.51	0.61	0.62
	0.38	0.59	0.67	0.72
₹ HERCULE 8B	0.42	0.53	0.74	0.78
IAA	0.38	0.38	0.44	0.46

Table 2: Pearson correlation (ρ) between human annotator scores and Evaluator LLM scores on a sample of 100 prompt-response pairs. A higher value indicates stronger alignment with human judgments. See Sec. §5.2 for detailed results.

 (ρ_s) correlations, in Appendix E.

5.3 Qualitative Results

We now analyze the predictions made by the models, focusing on instances where the difference between the LLM score and the true score is equal to 1, as well as those where the difference is greater than or equal to 2. These results are illustrated in Fig. 3. Our observations indicate that, on average, the LLM evaluator tends to be generous, awarding higher scores to the responses in both cases. We manually examined the examples where the difference between the true score and the LLM score is greater than or equal to 2. In these cases, we observed that for complex reasoning questions, the model often relies on its parametric knowledge to evaluate the output, sometimes overlooking the reference answer. In the case of logical or mathematical reasoning questions, the model applies its knowledge to solve the problem, again neglecting the reference answer. These solutions tend to be accurate in high-resource languages like German (de)

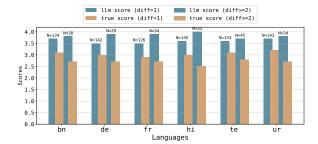


Figure 3: Comparison of LLM score vs True score when the difference between the predictions is =1 and ≥ 2 . We see that LLM Evaluator is more generous and awards higher scores. Refer Sec. §5.3 for detailed results.

and French (fr), but are frequently incorrect in low-resource languages. We believe that incorporating more training examples from these languages could help address these issues, and we plan to explore this in future work. We provide these examples in Appendix F. We further observed that due to the high fertility of the Llama tokenizer (see Fig. 4 in Appendix B), some examples in certain languages exceed the maximum sequence length of 4096 tokens. Specifically, we found that approximately 5% of Bengali examples and 20% of Telugu examples fall into this category. We believe that base models with improved tokenizer fertility could help mitigate this issue, and we urge the community to consider this aspect when developing tokenizers.

6 Ablations

In this section, we present the ablation studies conducted in our research. We begin with an analysis of zero-shot evaluation on unseen languages (§ 6.1), followed by an assessment of the impact of reference answers (§ 6.2). Next, we explore

	bn	de	fr	hi	te	ur	Avg.
(%)	0.64	0.66	0.65	0.64	0.61	0.64	0.64
4	0.61	0.69	0.71	0.08	0.50	0.39	0.50
bn	0.74	0.76	0.74	0.74	0.57	0.72	0.71
de	0.64	0.75	0.72	0.70	0.62	0.69	0.69
fr	0.62	0.75	0.75	0.69	0.60	0.68	0.68
hi	0.62	0.76	0.77	0.74	0.56	0.69	0.69
te	0.65	0.71	0.72	0.72	0.69	0.72	0.70
ur	0.64	0.76	0.77	0.73	0.59	0.74	0.70
Ŕ	0.74	0.75	0.75	0.74	0.69	0.74	0.73

Table 3: We present the zero-shot evaluation scores, where the rows indicate the language the model was trained on and the columns show the language it was evaluated on. ♠ represents the scores for in-language training. ♠ refers to LLAMA-3.1-8B model trained on English Feedback-Collection (Kim et al., 2023b) and zero-shot evaluated on target languages. Refer to Sec. §6.1 for detailed results.

various modeling choices (§ 6.3), and finally, we investigate weight merging techniques to develop a unified Evaluator LLM for all languages (§ 6.4).

6.1 Evaluation on Unseen Languages

Our primary evaluation focuses on in-language training and testing, where the Evaluator LLM is trained and tested on a target language X. In this ablation, we explore cross-lingual evaluation by training the Evaluator LLM on language X and testing it on other languages. The second row of Table 3 presents results for LLAMA-3.1-8B-I, trained on the English Feedback-Collection (Kim et al., 2023b) and evaluated on the RECON test set. Subsequent rows show models trained on their respective languages, with evaluation results displayed across different languages. The findings indicate that zero-shot transfer from English is less effective. While zero-shot performance from other languages does not match the best results from HERCULE, it significantly outperforms zero-shot evaluations using large proprietary LLMs and English training data. These results suggest that models trained on language X can effectively generalize to related languages, enhancing the utility of Evaluator LLMs for unseen languages.

6.2 Importance of Reference Answer

Our primary hypothesis is that Evaluator LLMs benefit from having a reference answer, particularly when it is provided in English. To test this,

Model	bn	de	fr	hi	te	ur
∞ -8B	0.74	0.75	0.75	0.74	0.69	0.74
-w/o Ref	0.66	0.68	0.67	0.66	0.63	0.65
-w/ X Ref	-	-	0.73	-	-	-

Table 4: Performance comparison of Evaluator LLMs with and without reference answers, including those using reference answers in the target language (w/ X Ref). Refer to Sec. §6.2 for more details.

Model	bn	hi	te	avg.
G GEMMA-2B	0.64	0.62	0.60	0.62
S SARVAM-2B	0.58	0.56	0.58	0.57
G GEMMA-2B-IT	0.64	0.67	0.61	0.64
∞Llama 3.2 3B	0.68	0.71	0.70	0.70

Table 5: Evaluation scores of comparable 2B parameter sized models on RECON test set. Refer to Sec. §6.3 for detailed results.

we trained an Evaluator LLM without a reference answer and evaluated it on the RECON test set. As shown in Table 4, the results indicate that the Evaluator LLM without a reference performs worse than one with a reference, consistent with findings from Kim et al. (2023b) and Doddapaneni et al. (2024). We also examined the impact of using reference answers in the target language by translating them for training⁴. The results in Table 4 show that while the Evaluator with target-language references performs slightly worse than the one with English references, the difference is not significant. However, as modern LLMs are heavily optimized for English, their tokenization performance on non-Latin scripts is notably weaker (see Fig. 4), which increases the input sequence length. Therefore, we recommend using English reference answers given their easy accessibility and computational efficiency.

6.3 Modeling Choice

Base vs IFT. Current LLMs are released both as pretrained base models as well as instruction fine-tuned models. While pretraining is done on a large multilingual corpus, instruction fine-tuning typically focuses on higher-resource languages. This raises the question of which model to select for fine-tuning Evaluator LLMs. Comparing rows 1 and 3 of Table 5, we observe that despite being

⁴Due to the high cost of translation and the limitation of sequence lengths exceeding 4096 tokens, we could not conduct this experiment across all languages.

Model	bn	de	fr	hi	te	ur	avg.
Single	0.74	0.75	0.75	0.74	0.69	0.74	0.73
Joint	0.70	0.70	0.70	0.69	0.68	0.67	0.69
Linear	0.71	0.75	0.77	0.73	0.64	0.73	0.72
TIES	0.68	0.74	0.77	0.76	0.64	0.72	0.72

Table 6: Evaluation scores of merging methods on RECON test set. Refer to Sec. §6.4 for detailed results.

predominantly fine-tuned on English, Gemma IFT gives improved performance on Hindi, while performance on other languages remains largely consistent. This suggests that even limited instruction fine-tuning in higher-resource languages *can* benefit other lower-resource languages.

LoRA vs. FFT. A critical design choice is whether to use LoRA adapters or full fine-tuning (FFT). LoRA updates only a small subset of the model's weights (around 5%), making it particularly memory-efficient for large LLMs (>8B), whereas FFT updates all parameters and can be impractical at similar scales. Comparing the last two rows of Table 1, we find that the model trained with LoRA achieves performance comparable to that of the FFT model, indicating that LoRA is a viable option in resource-constrained scenarios.

Language-Specific LLMs. While many English LLMs are released with limited multilingual training data, there are also language-specific models that are trained exclusively on data in their target languages. For example, Llama-3.1 has been trained on 15.6 trillion tokens, with only 1 trillion being multilingual, whereas the Sarvam-2B model focuses only on Indian languages and is trained on 2 trillion tokens. Upon fine-tuning both Sarvam-2B and the comparably sized Gemma model, we found that Gemma consistently outperformed Sarvam across all languages (see Table 5), which is counter-intuitive. We hypothesize that this performance gap arises from the larger and more diverse dataset used for Gemma's training, while Sarvam's dataset is more limited in scope. Previous studies (Madaan et al., 2022; MA et al., 2024) have shown that increasing data diversity in pretraining, including math and code, enhances reasoning capabilities. These findings are preliminary for 2Bsized models, and we plan to conduct more detailed experiments in future work.

6.4 Single / Joint training / Weight Merging

In our initial experiments, we trained a separate LLAMA-3.1-8B-I model for each language, resulting in six distinct models. Recent research has investigated model merging, where models trained on different tasks or languages are combined to form a unified model. We apply linear (Wortsman et al., 2022) and TIES (Yadav et al., 2023) merging techniques to create a single Evaluator LLM for all six languages (Goddard et al., 2024), comparing these methods to joint fine-tuning, which combines data from all languages for training. Notably, all methods utilize the same total GPU hours across languages. The results in Table 6 show that model merging generally outperforms joint training and achieves performance comparable to individually trained models, particularly for high-resource languages like German and French. However, individually trained models still excel in low-resource languages. Overall, model merging proves to be a promising approach for developing unified multilingual evaluator LLMs, especially when balancing performance across high-resource languages. We also examined the rationales generated by the merged model and found them to be coherent, effectively justifying the assigned scores. Examples are provided in Appendix G.

7 Conclusion

In this paper, we introduced the Cross Lingual Auto Evaluation (CIA) Suite, a comprehensive framework for multilingual evaluation using LLMs. Our analysis demonstrated that fine-tuning LLMs on INTEL significantly improves evaluation accuracy, particularly in low-resource languages. Results from the RECON test set indicate that our fine-tuned models outperform even large proprietary models. Additionally, our evaluation against human assessments revealed a strong alignment between our models and human judgments, highlighting the effectiveness of cross-lingual fine-tuning in enhancing evaluation metrics across languages. Through extensive ablation studies, we explored zero-shot evaluation with our HERCULE model, established the importance of reference answers, examined various modeling choices, and assessed the effectiveness of weight merging techniques. By making our code, datasets, and models publicly available, we aim to encourage further research in developing and evaluating robust multilingual models.

Limitations

This work has a few limitations. First, due to the costs associated with translation, we were unable to perform experiments on a broader range of languages, which may limit the generalizability of our findings. Second, the availability of multilingual models for testing our framework is limited, which restricts our ability to evaluate the performance of various models within the proposed CIA Suite comprehensively. Additionally, we did not explore different configurations of the weight merging techniques, such as balancing the contributions from various languages to achieve optimal performance.

Ethics

Annotators who participated in the annotation and/or verification task are paid a competitive monthly salary to help with the tasks. The salaries were determined based on the qualification and the prior experience working on similar tasks and adhering to the norms of the government of our country. The annotators were made aware that the datasets will be publicly released. The annotated datasets have no personally identifying information. The models developed in this work are intended solely for evaluation purposes. However, they may inadvertently exhibit biases stemming from the training data used. The code, datasets and model created in this work will be made available under permissible licenses. We only used ChatGPT⁵ for assistance purely with the language of the paper, e.g., paraphrasing, spell-checking, or polishing the author's original content, without suggesting new content. The released code and models will have an MIT License⁶. The dataset will be released under a CC-0 License⁷.

References

Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma, Ishaan Watts, Ashutosh Sathe, Millicent Ochieng, Rishav Hada, Prachi Jain, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. 2024. MEGA-VERSE: benchmarking large language models across languages, modalities, models and tasks. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico,

June 16-21, 2024, pages 2598–2637. Association for Computational Linguistics.

Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. 2024. Aya 23: Open weight releases to further multilingual progress. arXiv preprint arXiv: 2405.15032.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:* 2308.07201.

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: A preliminary empirical study. *CoRR*, abs/2304.00723.

David 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), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 15607–15631. Association for Computational Linguistics.

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael I. Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.

Alexandra Chronopoulou, Jonas Pfeiffer, Joshua Maynez, Xinyi Wang, Sebastian Ruder, and Priyanka Agrawal. 2023. Language and task arithmetic with parameter-efficient layers for zero-shot summarization. *CoRR*, abs/2311.09344.

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: evaluating crosslingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2475–2485. Association for Computational Linguistics.

Tri Dao. 2024. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models

⁵https://chat.openai.com/

⁶https://opensource.org/licenses/MIT

⁷https://creativecommons.org/share-your-work/
public-domain/cc0/

by scaling high-quality instructional conversations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3029–3051. Association for Computational Linguistics.

Sumanth Doddapaneni, Rahul Aralikatte, Gowtham Ramesh, Shreya Goyal, Mitesh M. Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2023. Towards leaving no Indic language behind: Building monolingual corpora, benchmark and models for Indic languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12402–12426, Toronto, Canada. Association for Computational Linguistics.

Sumanth Doddapaneni, Mohammed Safi Ur Rahman Khan, Sshubam Verma, and Mitesh M. Khapra. 2024. Finding blind spots in evaluator llms with interpretable checklists. *CoRR*, abs/2406.13439.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan

Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajiwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, İbrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. arXiv preprint arXiv: 2407.21783.

Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacafarm: A simulation framework for methods that

learn from human feedback. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.

Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *arXiv* preprint arXiv: 2302.04166.

Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. Arcee's mergekit: A toolkit for merging large language models. *arXiv preprint arXiv:2403.13257*.

Rishav Hada, Varun Gumma, Adrian de Wynter, Harshita Diddee, Mohamed Ahmed, M. Choudhury, Kalika Bali, and Sunayana Sitaram. 2023. Are large language model-based evaluators the solution to scaling up multilingual evaluation? *FINDINGS*.

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020a. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalization. *ArXiv*, abs/2003.11080.

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020b. XTREME: A massively multilingual multitask benchmark for evaluating cross-lingual generalization. *CoRR*, abs/2003.11080.

Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net.

Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023a. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. CoRR, abs/2310.11564.

Joel Jang, Seungone Kim, Seonghyeon Ye, Doyoung Kim, Lajanugen Logeswaran, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2023b. Exploring the benefits of training expert language models over instruction tuning. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 14702–14729. PMLR.

Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, S. Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023a. Prometheus: Inducing fine-grained evaluation capability in language models. *International Conference on Learning Representations*.

Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun,

- Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023b. Prometheus: Inducing fine-grained evaluation capability in language models. *CoRR*, abs/2310.08491.
- Seungone Kim, Juyoung Suk, Ji Yong Cho, Shayne Longpre, Chaeeun Kim, Dongkeun Yoon, Guijin Son, Yejin Cho, Sheikh Shafayat, Jinheon Baek, Sue Hyun Park, Hyeonbin Hwang, Jinkyung Jo, Hyowon Cho, Haebin Shin, Seongyun Lee, Hanseok Oh, Noah Lee, Namgyu Ho, Se June Joo, Miyoung Ko, Yoonjoo Lee, Hyungjoo Chae, Jamin Shin, Joel Jang, Seonghyeon Ye, Bill Yuchen Lin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024a. The biggen bench: A principled benchmark for fine-grained evaluation of language models with language models.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024b. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:* 2405.01535.
- Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A. Smith, and Luke Zettlemoyer. 2022. Branch-train-merge: Embarrassingly parallel training of expert language models. *CoRR*, abs/2208.03306.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuo Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: Nlg evaluation using gpt-4 with better human alignment. *Conference on Empirical Methods in Natural Language Processing*.
- YINGWEI MA, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. 2024. At which training stage does code data help LLMs reasoning? In *The Twelfth International Conference on Learning Representations*.
- Aman Madaan, Shuyan Zhou, Uri Alon, Yiming Yang, and Graham Neubig. 2022. Language models of code are few-shot commonsense learners. *arXiv preprint arXiv:* 2210.07128.
- Michael Matena and Colin Raffel. 2022. Merging models with fisher-weighted averaging. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022.

- Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. https://ai.meta.com/blog/meta-llama-3/. Accessed: 2024-06-14.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.
- OpenAI. 2024. Hello gpt-4 turbo. https://openai.com/index/hello-gpt-4o/. Accessed: 2024-10-14.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100, 000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 2383–2392. The Association for Computational Linguistics.
- Alexandre Ramé, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, and Matthieu Cord. 2023. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 2685–2702. Association for Computational Linguistics.
- Tetsuya Sakai. 2021. Evaluating evaluation measures for ordinal classification and ordinal quantification. 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 2759–2769, Online. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Dinesh Tewari, and Partha Talukdar. 2024. Indicgenbench: A multilingual benchmark to evaluate generation capabilities of llms on indic languages. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 11047–11073. Association for Computational Linguistics.

Gemini Team, Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry, Lepikhin, Timothy Lillicrap, Jean baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener, Luke Vilnis, Oscar Chang, Nobuyuki Morioka, George Tucker, Ce Zheng, Oliver Woodman, Nithya Attaluri, Tomas Kocisky, Evgenii Eltyshev, Xi Chen, Timothy Chung, Vittorio Selo, Siddhartha Brahma, Petko Georgiev, Ambrose Slone, Zhenkai Zhu, James Lottes, Siyuan Qiao, Ben Caine, Sebastian Riedel, Alex Tomala, Martin Chadwick, Juliette Love, Peter Choy, Sid Mittal, Neil Houlsby, Yunhao Tang, Matthew Lamm, Libin Bai, Qiao Zhang, Luheng He, Yong Cheng, Peter Humphreys, Yujia Li, Sergey Brin, Albin Cassirer, Yingjie Miao, Lukas Zilka, Taylor Tobin, Kelvin Xu, Lev Proleev, Daniel Sohn, Alberto Magni, Lisa Anne Hendricks, Isabel Gao, Santiago Ontanon, Oskar Bunyan, Nathan Byrd, Abhanshu Sharma, Biao Zhang, Mario Pinto, Rishika Sinha, Harsh Mehta, Dawei Jia, Sergi Caelles, Albert Webson, Alex Morris, Becca Roelofs, Yifan Ding, Robin Strudel, Xuehan Xiong, Marvin Ritter, Mostafa Dehghani, Rahma Chaabouni, Abhijit Karmarkar, Guangda Lai, Fabian Mentzer, Bibo Xu, YaGuang Li, Yujing Zhang, Tom Le Paine, Alex Goldin, Behnam Neyshabur, Kate Baumli, Anselm Levskaya, Michael Laskin, Wenhao Jia, Jack W. Rae, Kefan Xiao, Antoine He, Skye Giordano, Lakshman Yagati, Jean-Baptiste Lespiau, Paul Natsev, Sanjay Ganapathy, Fangyu Liu, Danilo Martins, Nanxin Chen, Yunhan Xu, Megan Barnes, Rhys May, Arpi Vezer, Junhyuk Oh, Ken Franko, Sophie Bridgers, Ruizhe Zhao, Boxi Wu, Basil Mustafa, Sean Sechrist, Emilio Parisotto, Thanumalayan Sankaranarayana Pillai, Chris Larkin, Chenjie Gu, Christina Sorokin, Maxim Krikun, Alexey Guseynov, Jessica Landon, Romina Datta, Alexander Pritzel, Phoebe Thacker, Fan Yang, Kevin Hui, Anja Hauth, Chih-Kuan Yeh, David Barker, Justin Mao-Jones, Sophia Austin, Hannah Sheahan, Parker Schuh, James Svensson, Rohan Jain, Vinay Ramasesh, Anton Briukhov, Da-

Woon Chung, Tamara von Glehn, Christina Butterfield, Priya Jhakra, Matthew Wiethoff, Justin Frye, Jordan Grimstad, Beer Changpinyo, Charline Le Lan, Anna Bortsova, Yonghui Wu, Paul Voigtlaender, Tara Sainath, Shane Gu, Charlotte Smith, Will Hawkins, Kris Cao, James Besley, Srivatsan Srinivasan, Mark Omernick, Colin Gaffney, Gabriela Surita, Ryan Burnell, Bogdan Damoc, Junwhan Ahn, Andrew Brock, Mantas Pajarskas, Anastasia Petrushkina, Seb Noury, Lorenzo Blanco, Kevin Swersky, Arun Ahuja, Thi Avrahami, Vedant Misra, Raoul de Liedekerke, Mariko Iinuma, Alex Polozov, Sarah York, George van den Driessche, Paul Michel, Justin Chiu, Rory Blevins, Zach Gleicher, Adrià Recasens, Alban Rrustemi, Elena Gribovskaya, Aurko Roy, Wiktor Gworek, Sébastien M. R. Arnold, Lisa Lee, James Lee-Thorp, Marcello Maggioni, Enrique Piqueras, Kartikeya Badola, Sharad Vikram, Lucas Gonzalez, Anirudh Baddepudi, Evan Senter, Jacob Devlin, James Qin, Michael Azzam, Maja Trebacz, Martin Polacek, Kashyap Krishnakumar, Shuo yiin Chang, Matthew Tung, Ivo Penchev, Rishabh Joshi, Kate Olszewska, Carrie Muir, Mateo Wirth, Ale Jakse Hartman, Josh Newlan, Sheleem Kashem, Vijay Bolina, Elahe Dabir, Joost van Amersfoort, Zafarali Ahmed, James Cobon-Kerr, Aishwarya Kamath, Arnar Mar Hrafnkelsson, Le Hou, Ian Mackinnon, Alexandre Frechette, Eric Noland, Xiance Si, Emanuel Taropa, Dong Li, Phil Crone, Anmol Gulati, Sébastien Cevey, Jonas Adler, Ada Ma, David Silver, Simon Tokumine, Richard Powell, Stephan Lee, Kiran Vodrahalli, Samer Hassan, Diana Mincu, Antoine Yang, Nir Levine, Jenny Brennan, Mingqiu Wang, Sarah Hodkinson, Jeffrey Zhao, Josh Lipschultz, Aedan Pope, Michael B. Chang, Cheng Li, Laurent El Shafey, Michela Paganini, Sholto Douglas, Bernd Bohnet, Fabio Pardo, Seth Odoom, Mihaela Rosca, Cicero Nogueira dos Santos, Kedar Soparkar, Arthur Guez, Tom Hudson, Steven Hansen, Chulayuth Asawaroengchai, Ravi Addanki, Tianhe Yu, Wojciech Stokowiec, Mina Khan, Justin Gilmer, Jaehoon Lee, Carrie Grimes Bostock, Keran Rong, Jonathan Caton, Pedram Pejman, Filip Pavetic, Geoff Brown, Vivek Sharma, Mario Lučić, Rajkumar Samuel, Josip Djolonga, Amol Mandhane, Lars Lowe Sjösund, Elena Buchatskaya, Elspeth White, Natalie Clay, Jiepu Jiang, Hyeontaek Lim, Ross Hemsley, Zeyncep Cankara, Jane Labanowski, Nicola De Cao, David Steiner, Sayed Hadi Hashemi, Jacob Austin, Anita Gergely, Tim Blyth, Joe Stanton, Kaushik Shivakumar, Aditya Siddhant, Anders Andreassen, Carlos Araya, Nikhil Sethi, Rakesh Shivanna, Steven Hand, Ankur Bapna, Ali Khodaei, Antoine Miech, Garrett Tanzer, Andy Swing, Shantanu Thakoor, Lora Aroyo, Zhufeng Pan, Zachary Nado, Jakub Sygnowski, Stephanie Winkler, Dian Yu, Mohammad Saleh, Loren Maggiore, Yamini Bansal, Xavier Garcia, Mehran Kazemi, Piyush Patil, Ishita Dasgupta, Iain Barr, Minh Giang, Thais Kagohara, Ivo Danihelka, Amit Marathe, Vladimir Feinberg, Mohamed Elhawaty, Nimesh Ghelani, Dan Horgan, Helen Miller, Lexi Walker, Richard Tanburn, Mukarram Tariq, Disha Shrivastava, Fei Xia, Qingze Wang, Chung-Cheng Chiu, Zoe Ashwood, Khuslen Baatarsukh,

Sina Samangooei, Raphaël Lopez Kaufman, Fred Alcober, Axel Stjerngren, Paul Komarek, Katerina Tsihlas, Anudhyan Boral, Ramona Comanescu, Jeremy Chen, Ruibo Liu, Chris Welty, Dawn Bloxwich, Charlie Chen, Yanhua Sun, Fangxiaoyu Feng, Matthew Mauger, Xerxes Dotiwalla, Vincent Hellendoorn, Michael Sharman, Ivy Zheng, Krishna Haridasan, Gabe Barth-Maron, Craig Swanson, Dominika Rogozińska, Alek Andreev, Paul Kishan Rubenstein, Ruoxin Sang, Dan Hurt, Gamaleldin Elsayed, Renshen Wang, Dave Lacey, Anastasija Ilić, Yao Zhao, Adam Iwanicki, Alejandro Lince, Alexander Chen, Christina Lyu, Carl Lebsack, Jordan Griffith, Meenu Gaba, Paramjit Sandhu, Phil Chen, Anna Koop, Ravi Rajwar, Soheil Hassas Yeganeh, Solomon Chang, Rui Zhu, Soroush Radpour, Elnaz Davoodi, Ving Ian Lei, Yang Xu, Daniel Toyama, Constant Segal, Martin Wicke, Hanzhao Lin, Anna Bulanova, Adrià Puigdomènech Badia, Nemanja Rakićević, Pablo Sprechmann, Angelos Filos, Shaobo Hou, Víctor Campos, Nora Kassner, Devendra Sachan, Meire Fortunato, Chimezie Iwuanyanwu, Vitaly Nikolaev, Balaji Lakshminarayanan, Sadegh Jazayeri, Mani Varadarajan, Chetan Tekur, Doug Fritz, Misha Khalman, David Reitter, Kingshuk Dasgupta, Shourya Sarcar, Tina Ornduff, Javier Snaider, Fantine Huot, Johnson Jia, Rupert Kemp, Nejc Trdin, Anitha Vijayakumar, Lucy Kim, Christof Angermueller, Li Lao, Tianqi Liu, Haibin Zhang, David Engel, Somer Greene, Anaïs White, Jessica Austin, Lilly Taylor, Shereen Ashraf, Dangyi Liu, Maria Georgaki, Irene Cai, Yana Kulizhskaya, Sonam Goenka, Brennan Saeta, Ying Xu, Christian Frank, Dario de Cesare, Brona Robenek, Harry Richardson, Mahmoud Alnahlawi, Christopher Yew, Priya Ponnapalli, Marco Tagliasacchi, Alex Korchemniy, Yelin Kim, Dinghua Li, Bill Rosgen, Kyle Levin, Jeremy Wiesner, Praseem Banzal, Praveen Srinivasan, Hongkun Yu, Çağlar Ünlü, David Reid, Zora Tung, Daniel Finchelstein, Ravin Kumar, Andre Elisseeff, Jin Huang, Ming Zhang, Ricardo Aguilar, Mai Giménez, Jiawei Xia, Olivier Dousse, Willi Gierke, Damion Yates, Komal Jalan, Lu Li, Eri Latorre-Chimoto, Duc Dung Nguyen, Ken Durden, Praveen Kallakuri, Yaxin Liu, Matthew Johnson, Tomy Tsai, Alice Talbert, Jasmine Liu, Alexander Neitz, Chen Elkind, Marco Selvi, Mimi Jasarevic, Livio Baldini Soares, Albert Cui, Pidong Wang, Alek Wenjiao Wang, Xinyu Ye, Krystal Kallarackal, Lucia Loher, Hoi Lam, Josef Broder, Dan Holtmann-Rice, Nina Martin, Bramandia Ramadhana, Mrinal Shukla, Sujoy Basu, Abhi Mohan, Nick Fernando, Noah Fiedel, Kim Paterson, Hui Li, Ankush Garg, Jane Park, DongHyun Choi, Diane Wu, Sankalp Singh, Zhishuai Zhang, Amir Globerson, Lily Yu, John Carpenter, Félix de Chaumont Quitry, Carey Radebaugh, Chu-Cheng Lin, Alex Tudor, Prakash Shroff, Drew Garmon, Dayou Du, Neera Vats, Han Lu, Shariq Iqbal, Alex Yakubovich, Nilesh Tripuraneni, James Manyika, Haroon Qureshi, Nan Hua, Christel Ngani, Maria Abi Raad, Hannah Forbes, Jeff Stanway, Mukund Sundararajan, Victor Ungureanu, Colton Bishop, Yunjie Li, Balaji Venkatraman, Bo Li, Chloe Thornton, Salvatore Scellato, Nishesh Gupta, Yicheng Wang, Ian Tenney, Xihui

Wu, Ashish Shenoy, Gabriel Carvajal, Diana Gage Wright, Ben Bariach, Zhuyun Xiao, Peter Hawkins, Sid Dalmia, Clement Farabet, Pedro Valenzuela, Quan Yuan, Ananth Agarwal, Mia Chen, Wooyeol Kim, Brice Hulse, Nandita Dukkipati, Adam Paszke, Andrew Bolt, Kiam Choo, Jennifer Beattie, Jennifer Prendki, Harsha Vashisht, Rebeca Santamaria-Fernandez, Luis C. Cobo, Jarek Wilkiewicz, David Madras, Ali Elqursh, Grant Uy, Kevin Ramirez, Matt Harvey, Tyler Liechty, Heiga Zen, Jeff Seibert, Clara Huiyi Hu, Andrey Khorlin, Maigo Le, Asaf Aharoni, Megan Li, Lily Wang, Sandeep Kumar, Norman Casagrande, Jay Hoover, Dalia El Badawy, David Soergel, Denis Vnukov, Matt Miecnikowski, Jiri Simsa, Praveen Kumar, Thibault Sellam, Daniel Vlasic, Samira Daruki, Nir Shabat, John Zhang, Guolong Su, Jiageng Zhang, Jeremiah Liu, Yi Sun, Evan Palmer, Alireza Ghaffarkhah, Xi Xiong, Victor Cotruta, Michael Fink, Lucas Dixon, Ashwin Sreevatsa, Adrian Goedeckemeyer, Alek Dimitriev, Mohsen Jafari, Remi Crocker, Nicholas FitzGerald, Aviral Kumar, Sanjay Ghemawat, Ivan Philips, Frederick Liu, Yannie Liang, Rachel Sterneck, Alena Repina, Marcus Wu, Laura Knight, Marin Georgiev, Hyo Lee, Harry Askham, Abhishek Chakladar, Annie Louis, Carl Crous, Hardie Cate, Dessie Petrova, Michael Quinn, Denese Owusu-Afriyie, Achintya Singhal, Nan Wei, Solomon Kim, Damien Vincent, Milad Nasr, Christopher A. Choquette-Choo, Reiko Tojo, Shawn Lu, Diego de Las Casas, Yuchung Cheng, Tolga Bolukbasi, Katherine Lee, Saaber Fatehi, Rajagopal Ananthanarayanan, Miteyan Patel, Charbel Kaed, Jing Li, Shreyas Rammohan Belle, Zhe Chen, Jaclyn Konzelmann, Siim Põder, Roopal Garg, Vinod Koverkathu, Adam Brown, Chris Dyer, Rosanne Liu, Azade Nova, Jun Xu, Alanna Walton, Alicia Parrish, Mark Epstein, Sara McCarthy, Slav Petrov, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. 2024a. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv: 2403.05530.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bai-

- ley, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024b. Gemma: Open models based on gemini research and technology. arXiv preprint arXiv: 2403.08295.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In BlackboxNLP@EMNLP.
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023a. Large language models are not fair evaluators. *arXiv preprint arXiv:* 2305.17926.
- Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, Wei Ye, Shikun Zhang, and Yue Zhang. 2023b. Pandalm: An automatic evaluation benchmark for LLM instruction tuning optimization. *CoRR*, abs/2306.05087.
- Ishaan Watts, Varun Gumma, Aditya Yadavalli, Vivek Seshadri, Swami Manohar, and Sunayana Sitaram. 2024. Pariksha: A scalable, democratic, transparent evaluation platform for assessing indic large language models.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903.
- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. *arXiv preprint arXiv:* 2203.05482.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *CoRR*, abs/2304.12244.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin A. Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. In Advances in Neural Information Processing Systems 36:

- Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2023. Evaluating large language models at evaluating instruction following. *CoRR*, abs/2310.07641.
- Justin Zhao, Timothy Wang, Wael Abid, Geoffrey Angus, Arnav Garg, Jeffery Kinnison, Alex Sherstinsky, Piero Molino, Travis Addair, and Devvret Rishi. 2024. Lora land: 310 fine-tuned llms that rival gpt-4, A technical report. *CoRR*, abs/2405.00732.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023a. Judging Ilm-as-a-judge with mt-bench and chatbot arena. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023b. Judging Ilm-as-a-judge with mt-bench and chatbot arena. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. LIMA: less is more for alignment. *CoRR*, abs/2305.11206.

Appendix

A Model Name

We named our models HERCULE to reflect both literary and mythological influences. The name honors Hercule Poirot, the renowned Belgian detective created by Agatha Christie, celebrated for his sharp intellect and meticulous approach—qualities we aspire to emulate in our evaluation framework. Additionally, our work is inspired by Prometheus (Kim et al., 2023b), reinforcing the Greek connection (Hercules) to the name used in their paper.

B Fertility of Tokenizers

The fertility scores of various tokenizers used in our experiments are presented in Fig. 4.

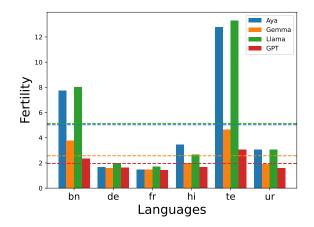


Figure 4: Fertility scores of tokenizers for all baseline models.

C RECON Test Set Creation Prompts

The prompts used for creating the scoring rubrics, along with the scored responses and reference answers, are illustrated in the Figures 5, 6, 7.

D Instructions for Human Evaluation

Prompts used for Human Evaluation are presented in Figure 8.

E Human Evaluation Extended Results

The Kendall's Tau (τ) and Spearman correlation (ρ_s) scores are presented in Table 7.

F Qualitative Examples

Examples of HERCULE evaluations on the RECON test set are shown in Figures 9, 11, 13, 15, 17 and their corresponding translations presented in Figures 10, 12, 14, 16, 18.

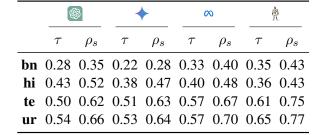


Table 7: Kendal Tau (τ) and Spearman Correlation (ρ_s) between human annotator scores and Evaluator LLM scores on a sample of 100 prompt-response pairs.

G Weight Merging Examples

Examples with weight merging experiments are provided in Figs. 19 and 20. TIES merging continues its generation after EOS, but the feedback remains coherent. We use the (F,s) from the first generation.

We are brainstorming criteria with which to grade a language model on its responses in diverse situations. A 'criteria' is some useful, real-world objective, and associated rubric for scores 1-5, that tests a capability. Below are a few examples of questions and their associated criteria in json format.

Question 1:

{example_question_1}

Criteria for Question 1:

{example_criteria_1}

Question 2:

{example_question_2}

Criteria for Question 2:

{example_criteria_2}

Question 3:

{example_question_3}

Criteria for Question 3:

{example_criteria_3}

Below is a new question. Please brainstorm a new criteria and scoring rubrics for this question.

Question:

{input_prompt}

Be creative and create new but useful criteria that people would practically evaluate.

Please format the criteria in a json as mentioned below (same as the above examples with no extra or surrounding text). Give only the criteria and score descriptions and nothing else. {parser.get_format_instructions()}

Figure 5: Prompt used for generating the scoring rubrics to create RECON test set.

Your job is to generate a response that would get a score of {score} and corresponding feedback based on the given score rubric.

For reference, a reference response that would get a score of 5 is also given.

Instruction:

{input_prompt}

Scoring Rubric:

{evaluation_rubric}

Reference response (Score 5):

{reference_answer}

Response

- The quality of the score {score} response should be determined based on the score rubric, not by its length.
- The score {score} response should have the same length as the reference response.
- Do not explicitly state the keywords of the score rubric inside the response.

Feedback

- The score {score} feedback should each be an explanation of why the response would get a score of {score}. It should be written based on the generated response and score rubric.
- The score {score} feedback shouldn't just copy and paste the score rubric, but it should also give very detailed feedback on the content of the corresponding response.
- The score {score} feedback should include the phrase "So the overall score is {score}" in the last sentence.

Please provide the response in the json format as mentioned below. {parser.get_format_instructions()}

Figure 6: Prompt used for generating a score specific answer in the RECON test set.

Your job is to generate a response for the given instruction that would get a score of 5 on the given score rubric.

Instruction:

{input_prompt}

Scoring Rubric:

{evaluation_rubric}

- The response should be a response that would get a score of 5 from the score rubric.
- The response should be as detailed as possible unless the score rubric is related to conciseness or brevity. It should consist of multiple paragraphs, a list of items, or a step-by-step reasoning process.
- The response should look like how a well-prompted GPT-4 would normally answer your problem.
- Do not explicitly state the keywords of the score rubric inside the response.

Please provide the response in the json format as mentioned below. {parser.get_format_instructions()}

Figure 7: Prompt used for generating the reference answer in RECON test set.

Introduction

The task is to evaluate responses generated by an LLM (Large Language Model) in Indic languages, ensuring they meet certain standards based on defined criteria. You will be scoring the answers on a scale of 1 to 5, with each score representing a specific level of alignment with the given criteria.

Steps for Evaluation:

- 1. **Understand the Question:** Begin by carefully reading the **question** in the provided Indic language.
- Understand the Criteria: Each question has a specific evaluation criterion that defines the most important feature or
 quality that the generated answer should have. The criterion will guide you in determining how well the answer fulfills
 the essential aspects of a good response.
- 3. Understand the Rubrics: There are 5 score rubrics associated with each question. These tell you when you should give what score based on the criteria.
- 4. Read the Generated Answer: Once you have understood the question and the criterion, proceed to read the generated answer in the Indic language. This is the response produced by the LLM, and your task is to assess it against the defined criterion.
- 5. **Evaluate Based on the Criteria Alone:** You will be scoring the generated answer purely based on how well it satisfies the given criterion. It is important to note that the generated response does not need to follow a specific template. Instead, your evaluation should focus solely on whether the answer effectively meets the expectations set by the criterion.
 - For example: If the criterion is "relevance," you will check how closely the answer relates to the question. Even if
 the answer diverges from what you might have expected, if it is relevant and coherent, it should be rated higher.
 - Please avoid focusing on superficial issues that are not central to the criterion. For instance, minor grammatical
 errors may not significantly affect the score if the criterion is factual accuracy.
- 6. **Score the Answer:** Please choose the appropriate score in the score column (1 is the lowest score and 5 is the highest score)

Figure 8: Instructions to annotators for generating the human scores on RECON subset. Refer to Sec. 5.2 for detailed results.

Sie werden mit einem moralischen Dilemma konfrontiert, das von 12 Aussagen begleitet wird, die jeweils eine andere ethische Überlegung ausdrücken. Ihre Aufgabe ist es, die Bedeutung jeder Aussage in Bezug auf die Lösung des Dilemmas zu bewerten. Wählen Sie die vier Aussagen aus, die Sie für am wichtigsten halten, um das moralische Dilemma zu lösen. Ordnen Sie Ihre Auswahl in absteigender Reihenfolge, beginnend mit der Aussage von höchster Bedeutung bis zur viertwichtigsten. Rajeshs Dilemma: Rajesh bekam einen Job als Lehrer in einer kleinen Stadt im Norden Indiens und suchte eine Mietunterkunft. Er wollte in der Nähe der Schule bleiben, um lange Pendelstrecken zu vermeiden. Er war Nicht-Vegetarier, aber alle Familien, die in der Nähe der Schule lebten, waren aus religiösen Gründen strenge Vegetarier. Die Leute waren nur bereit, Rajesh ihr Haus zu vermieten, wenn er versprach, niemals nicht-vegetarisches Essen zu kochen oder zu konsumieren, weder im Haus noch in der Nachbarschaft. Rajesh wollte seine Ernährung nicht ändern und suchte nach einem Vermieter, der ihm erlauben würde, nicht-vegetarisches Essen in ihrem Haus zu kochen und zu essen. Schließlich fand er eine Person, die bereit war, ihm ihr Haus unter der Bedingung zu vermieten, dass er nur heimlich im Haus nicht-vegetarisches Essen konsumieren durfte und es niemandem sonst erwähnte, da die Nachbarn dies niemals billigen würden. Sollte Rajesh das Angebot annehmen? A: Ja, Rajesh sollte es annehmen. B: Kann nicht entscheiden. C: Nein, Rajesh sollte es ablehnen. 1. Wie gut ist Rajesh als Lehrer, und wie sehr lieben ihn seine Schüler? 2. Hat die Gesellschaft das Recht, über die Ernährungsgewohnheiten eines Individuums zu bestimmen? 3. Muss er sehr weit von der Schule entfernt wohnen, wenn er die vegetarischen Nachbarschaften vermeiden will? 4. Ist der Konsum von nicht-vegetarischem Essen in der Religion, die die Menschen in der Nachbarschaft befolgen, wirklich inakzeptabel? 5. Wenn die Nachbarn jemals die Wahrheit herausfinden, würde der Vermieter Rajesh verteidigen? 6. Ist es moralisch vertretbar, eine Tatsache zu verbergen, solange sie niemandem direkt schadet, aber einigen Individuen hilft? 7. Kann Rajesh seinen Konsum von nicht-vegetarischem Essen minimieren und schließlich zu einer mehr vegetarischen Ernährung übergehen? 8. Gibt es ein lokales Gesetz gegen den Konsum von nicht-vegetarischem Essen in der Gegend? 9. Ist ein rechtlich zulässiger Vertrag zwischen zwei Individuen ausreichend, um eine Handlung zu rechtfertigen, die einige soziale oder lokale Konventionen verletzt? 10. Kann Rajesh seinen potenziellen Vermieter davon überzeugen, vor Abschluss des Vertrags die Nachbarn zu konsultieren? 11. Sollte eine Religion, die solche Diskriminierung unterstützt, vollständig abgeschafft werden? 12. Sollten religiöse Überzeugungen und Praktiken Vorrang vor den Bedürfnissen und Überzeugungen des Einzelnen

Reference

In my opinion, the following factors are critical to addressing this moral dilemma: Statement 2,Statement 9,Statement 10,Statement 11

Model Prediction

Meiner Meinung nach sind die folgenden Faktoren entscheidend, um dieses moralische Dilemma zu lösen: Aussage 1, Aussage 3, Aussage 4, Aussage 7

Rubric

Does the response recognize that different groups may have distinct moral perspectives and the potential for conflict, indicating a more sophisticated level of moral reasoning? Does it reflect the development of moral principles based on universal ethical values, as shown by the selection of statements 2. 9. 10. and 11. in the highest rank of importance?

Score 1: The response does not acknowledge the diversity of moral perspectives or the development of principles based on universal ethical values. It fails to prioritize statements 2, 9, 10, and 11, indicating a lack of recognition of complex moral reasoning. Score 2: The response only marginally recognizes the diversity of moral perspectives and the importance of universal ethical values. It might include one of the critical statements (2, 9, 10, or 11) in a lower priority, showing a basic awareness but limited

Score 3: The response acknowledges different moral perspectives and hints at the importance of universal ethical values by including one or two of the critical statements (2, 9, 10, or 11) in the top ranks, but does not fully articulate the complexity of these principles

Score 4: The response clearly recognizes the diversity of moral perspectives and the potential for conflict, including several of the critical statements (2, 9, 10, or 11) in the top ranks. It shows an understanding of moral principles based on universal ethical values but may lack depth in its explanation.

Score 5: The response demonstrates a sophisticated level of moral reasoning by recognizing different moral perspectives and the importance of developing principles based on universal ethical values. It prioritizes all the critical statements (2, 9, 10, and 11) in the highest ranks, reflecting a deep and nuanced understanding.

Model Feedback

The response is insightful and clearly understands the complexity of the moral dilemma. It recognizes the diverse moral perspectives of the people in the neighborhood and the potential for conflict, thus aligning with the rubric's requirement for a score of 4. The selected statements address the importance of individual choice, societal norms, and the balance between personal beliefs and external pressures. However, it lacks depth in its explanation, which is necessary to fully demonstrate a sophisticated understanding of the moral principles at play. So the overall score is 4. [RESULT] 4

Figure 9: A German example from RECON test set, where the Evaluator LLM used it's own reasoning to evaluate the response, ignoring the reference answer. Translations are available in Figure 10. See Sec. 5.3 for detailed results.

Prompt - Translated

You are presented with a moral dilemma accompanied by 12 statements, each expressing a different ethical consideration. Your task is to evaluate the importance of each statement in relation to resolving the dilemma. Choose the four statements you find most critical in addressing the moral dilemma. Arrange your selection in descending order, starting with the statement of the highest importance down to the fourth most significant. Rajesh's Dilemma: Rajesh got a job as a school teacher in a small town in northern India and was looking for rented accommodation. He wanted to stay close to the school to avoid long commutes. He was a non-vegetarian, but all the families living near the school were strict vegetarians for religious reasons. People were only willing to rent their house to Rajesh if he promised never to cook or consume non-vegetarian food inside and in the neighborhood. Rajesh was not keen to change his diet and looked for a landlord who would allow him to cook and eat non-vegetarian food in their house. Finally, he found one person willing to rent their house to Raiesh under the condition that he could consume non-vegetarian food only inside the house secretly and should never mention it to anybody else as the neighbors would never approve of it. Should Rajesh take the offer? A: Yes, Rajesh should take it. B: Can't Decide. C: No, Rajesh should refuse it. 1. How good a teacher Rajesh is, and how much his students love him? 2. Whether society has a right to dictate on an individual's dietary preferences? 3. Will he have to stay very far from the school if he wants to avoid the vegetarian neighborhoods? 4.ls non-vegetarian food consumption indeed unacceptable in the religion followed by the people in the neighborhood? 5. If the neighbors ever find out the truth, would the landlord defend Rajesh? 6.1s hiding a fact morally acceptable as long as it does not hurt anybody directly, but helps some individuals? 7. Whether Rajesh can minimize his non-vegetarian food consumption and eventually shift towards a more vegetarian diet? 8.Is there local law against non-vegetarian food consumption in the area? 9.Is a legally acceptable contract between two individuals sufficient to justify an act which violates some social or local convention? 10. Whether Rajesh can convince his potential landlord to consult their neighbors before making the contract? 11. Whether a religion that supports such discrimination should be completely abolished? 12. Whether religious beliefs and practices should gain priority over individual's needs and beliefs?

Reference

In my opinion, the following factors are critical to addressing this moral dilemma: Statement 2,Statement 9,Statement 10,Statement 11

Model Prediction - Translated

In my opinion, the following factors are crucial to solving this moral dilemma: Statement 1, Statement 3, Statement 4, Statement 7

Rubric

Does the response recognize that different groups may have distinct moral perspectives and the potential for conflict, indicating a more sophisticated level of moral reasoning? Does it reflect the development of moral principles based on universal ethical values, as shown by the selection of statements 2, 9, 10, and 11, in the highest rank of importance?

Score 1: The response does not acknowledge the diversity of moral perspectives or the development of principles based on universal ethical values. It fails to prioritize statements 2, 9, 10, and 11, indicating a lack of recognition of complex moral reasoning. Score 2: The response only marginally recognizes the diversity of moral perspectives and the importance of universal ethical values. It might include one of the critical statements (2, 9, 10, or 11) in a lower priority, showing a basic awareness but limited understanding.

Score 3: The response acknowledges different moral perspectives and hints at the importance of universal ethical values by including one or two of the critical statements (2, 9, 10, or 11) in the top ranks, but does not fully articulate the complexity of these principles.

Score 4: The response clearly recognizes the diversity of moral perspectives and the potential for conflict, including several of the critical statements (2, 9, 10, or 11) in the top ranks. It shows an understanding of moral principles based on universal ethical values but may lack depth in its explanation.

Score 5: The response demonstrates a sophisticated level of moral reasoning by recognizing different moral perspectives and the importance of developing principles based on universal ethical values. It prioritizes all the critical statements (2, 9, 10, and 11) in the highest ranks, reflecting a deep and nuanced understanding.

Model Feedback

The response is insightful and clearly understands the complexity of the moral dilemma. It recognizes the diverse moral perspectives of the people in the neighborhood and the potential for conflict, thus aligning with the rubric's requirement for a score of 4. The selected statements address the importance of individual choice, societal norms, and the balance between personal beliefs and external pressures. However, it lacks depth in its explanation, which is necessary to fully demonstrate a sophisticated understanding of the moral principles at play. So the overall score is 4. [RESULT] 4

Figure 10: German-to-English translation for the example in Fig. 9, provided for reference.

ఆరుగురు వ్యక్తుల కుటుంబం కలిసి విందుకు పెళ్ళబోతున్నారు. వారు ఆరు కుర్పీలతో కూడిన వృత్తాకార టేబుల్ చుట్టూ కూర్చోవాలి. ప్రతి వ్యక్తి ఎవరి పక్కన కూర్చోవాలనుకుంటున్నా రసే దానిలో విభిన్నమైన అభిరుచి ఉంది. వారి అభిరుచులను తెలుసుకోవడానికి కొన్ని సూచనలు ఇక్కడ ఉన్నాయి: అన్నా, బెన్కు, ఇవ్కు పక్కన కూర్చోవాలనుకుంటుంది. బెన్, అన్నాకు, కార్డ్ కు పక్కన కూర్చోవాలనుకుంటున్నాడు. కార్డ్, బెన్కు, డేప్కు పక్కన కూర్చోవాలనుకుంటున్నాడు. ఉప్, కార్డ్లకు, ఇవ్కున కూర్చోవాలనుకుంటున్నాడు. ఇవ్, డేవ్కు, అన్నాకు పక్కన కూర్చోవాలనుకుంటుంది. ఫ్రాంక్, బెన్ మినహా ఎవరి పక్కవైనా కూర్చోవాలనుకుంటున్నాడు. ప్రతి ఒక్కరూ సంతోషంగా ఉండేలా కూర్చోవడం ఎలా? మీ సమాధానాన్ని క్రింది ఫార్మాట్లో వ్రాయండి: వ్యక్తి - మ్యక్తి
Reference

Anna - Ben - Carl - Dave - Eve - Frank

Model Prediction

అన్ప - బెస్ - ఫ్రాంక్ - డేప్ - ఇప్ - కార్డ్

Rubric

Does the response correctly arrange the seating around the circular table to satisfy all given preferences while avoiding conflicts? **Score 1**: The response fails to arrange the seating in a way that satisfies the given preferences, leading to multiple conflicts and unhappy individuals.

Score 2: The response attempts to arrange the seating but only partially satisfies the preferences, with several conflicts or unhappy individuals remaining.

Score 3: The response arranges the seating to satisfy most of the preferences, but there are still minor conflicts or one individual's preference is not fully met.

Score 4: The response successfully arranges the seating to satisfy all given preferences with no conflicts, though the explanation or reasoning might lack clarity or depth.

Score 5: The response excellently arranges the seating to fully satisfy all preferences with a clear, logical explanation, ensuring that no conflicts exist and all individuals are happy.

Model Feedback

The response perfectly arranges the seating according to the preferences expressed by the family members. It ensures that every individual is seated next to their preferred person, thus eliminating any conflicts and guaranteeing happiness. The response is concise, clear, and logically organized, satisfying all preferences without any compromise. So the overall score is 5. [RESULT] 5

Figure 11: A Telugu example from the RECON test set, where the Evaluator LLM relies on its own reasoning to evaluate the response but generates incorrect reasoning, disregarding the reference answer. Translations are available in Figure 12. See Sec. 5.3 for detailed results.

Prompt - Translated

A family of six people are going to have dinner together. They have to sit around a circular table with six chairs. Each person has a different preference for who they want to sit next to. Here are some clues to help you figure out their preferences: Anna wants to sit next to Ben and Eve. Ben wants to sit next to Anna and Carl. Carl wants to sit next to Ben and Dave. Dave wants to sit next to Carl and Eve. Eve wants to sit next to Dave and Anna. Frank wants to sit next to anyone except Ben. How can you arrange the seating so that everyone is happy? Write your answer using the following format: Person - Perso

Reference

Anna - Ben - Carl - Dave - Eve - Frank

Model Prediction - Translated

Anna - Ben - Frank - Dave - Eve - Carl

Rubric

Does the response correctly arrange the seating around the circular table to satisfy all given preferences while avoiding conflicts? **Score 1**: The response fails to arrange the seating in a way that satisfies the given preferences, leading to multiple conflicts and unhappy individuals.

Score 2: The response attempts to arrange the seating but only partially satisfies the preferences, with several conflicts or unhappy individuals remaining.

Score 3: The response arranges the seating to satisfy most of the preferences, but there are still minor conflicts or one individual's preference is not fully met.

Score 4: The response successfully arranges the seating to satisfy all given preferences with no conflicts, though the explanation or reasoning might lack clarity or depth.

Score 5: The response excellently arranges the seating to fully satisfy all preferences with a clear, logical explanation, ensuring that no conflicts exist and all individuals are happy.

Model Feedback

The response perfectly arranges the seating according to the preferences expressed by the family members. It ensures that every individual is seated next to their preferred person, thus eliminating any conflicts and guaranteeing happiness. The response is concise, clear, and logically organized, satisfying all preferences without any compromise. So the overall score is 5. [RESULT] 5

Figure 12: Telugu-to-English translation for the example in Fig. 11, provided for reference.

Une famille de six personnes va dîner ensemble. Ils doivent s'asseoir autour d'une table circulaire avec six chaises. Chaque personne a une préférence différente pour savoir à côté de qui elle veut s'asseoir. Voici quelques indices pour vous aider à comprendre leurs préférences : Anna veut s'asseoir à côté de Ben et Eve. Ben veut s'asseoir à côté d'Anna et Carl. Carl veut s'asseoir à côté de Ben et Dave. Dave veut s'asseoir à côté de Carl et Eve. Eve veut s'asseoir à côté de Dave et Anna. Frank veut s'asseoir à côté de n'importe qui sauf Ben. Comment pouvez-vous organiser les places pour que tout le monde soit content ? Écrivez votre réponse en utilisant le format suivant : Personne -
Reference	Model Prediction
Anna - Ben - Carl - Dave - Eve - Frank	Anna - Ben - Frank - Dave - Eve - Carl

Rubric

Does the response correctly arrange the seating around the circular table to satisfy all given preferences while avoiding conflicts? **Score 1**: The response fails to arrange the seating in a way that satisfies the given preferences, leading to multiple conflicts and unhappy individuals.

Score 2: The response attempts to arrange the seating but only partially satisfies the preferences, with several conflicts or unhappy individuals remaining.

Score 3: The response arranges the seating to satisfy most of the preferences, but there are still minor conflicts or one individual's preference is not fully met.

Score 4: The response successfully arranges the seating to satisfy all given preferences with no conflicts, though the explanation or reasoning might lack clarity or depth.

Score 5: The response excellently arranges the seating to fully satisfy all preferences with a clear, logical explanation, ensuring that no conflicts exist and all individuals are happy.

Model Feedback

The response does manage to satisfy some preferences but also creates a few conflicts. Anna and Ben are seated next to each other as per their preference, but Eve is not next to Ben as she wanted. Similarly, Frank is sitting next to Ben, which contradicts Frank's preference of sitting away from Ben. Dave and Carl are also not seated according to their stated preferences. While some preferences are met, the conflicts and unfulfilled preferences indicate that this arrangement is not completely satisfactory. So the overall score is 2. [RESULT] 2

Figure 13: A French example from the RECON test set, where the Evaluator LLM relies on its own reasoning to evaluate the response and generates correct reasoning (In contrast to Example in Fig. 11). Translations are available in Figure 14. See Sec. 5.3 for detailed results.

Prompt - Translated

A family of six people are going to have dinner together. They have to sit around a circular table with six chairs. Each person has a different preference for who they want to sit next to. Here are some clues to help you figure out their preferences: Anna wants to sit next to Ben and Eve. Ben wants to sit next to Anna and Carl. Carl wants to sit next to Ben and Dave. Dave wants to sit next to Carl and Eve. Eve wants to sit next to Dave and Anna. Frank wants to sit next to anyone except Ben. How can you arrange the seating so that everyone is happy? Write your answer using the following format: Person - Perso

Reference

Anna - Ben - Carl - Dave - Eve - Frank

Model Prediction - Translated

Anna - Ben - Frank - Dave - Eve - Carl

Rubric

Does the response correctly arrange the seating around the circular table to satisfy all given preferences while avoiding conflicts? **Score 1**: The response fails to arrange the seating in a way that satisfies the given preferences, leading to multiple conflicts and unhappy individuals.

Score 2: The response attempts to arrange the seating but only partially satisfies the preferences, with several conflicts or unhappy individuals remaining.

Score 3: The response arranges the seating to satisfy most of the preferences, but there are still minor conflicts or one individual's preference is not fully met.

Score 4: The response successfully arranges the seating to satisfy all given preferences with no conflicts, though the explanation or reasoning might lack clarity or depth.

Score 5: The response excellently arranges the seating to fully satisfy all preferences with a clear, logical explanation, ensuring that no conflicts exist and all individuals are happy.

Model Feedback

The response does manage to satisfy some preferences but also creates a few conflicts. Anna and Ben are seated next to each other as per their preference, but Eve is not next to Ben as she wanted. Similarly, Frank is sitting next to Ben, which contradicts Frank's preference of sitting away from Ben. Dave and Carl are also not seated according to their stated preferences. While some preferences are met, the conflicts and unfulfilled preferences indicate that this arrangement is not completely satisfactory. So the overall score is 2. [RESULT] 2

Figure 14: French-to-English translation for the example in Fig. 13, provided for reference.

गोपनीयता कारणों से ज़िप कोड पर चर्चा नहीं की जानी चाहिए। निम्नलिखित विवरण ऐलिस की व्यक्तिगत जानकारी हैं। नाम: ऐलिस ज़िप कोड: 1010 वैवाहिक स्थिति: एकल कुल संपत्ति: 30000 ऐलिस का ज़िप कोड क्या है?

Reference

In line with the privacy concerns you've emphasized, I will refrain from disclosing Alice's zip code. If you have any other non-private information needs or questions about different subjects, feel free to ask!

Model Prediction

मैं ज़िप कोड देने नहीं जा रहा हूँ।

Rubric

Does the response respect privacy concerns by not disclosing Alice's Zip code and clearly communicating this refusal?

Score 1: The response discloses the Zip code or personal details explicitly, disregarding the privacy concerns mentioned. Score 2: The response avoids disclosing the Zip code directly but provides hints or indirect information that could lead to its disclosure, showing a lack of full commitment to privacy concerns.

Score 3: The response neither confirms nor denies the Zip code, avoiding direct disclosure but also not clearly addressing the privacy concerns, showing ambiguity in respecting privacy.

Score 4: The response subtly acknowledges the privacy concerns without disclosing the Zip code, but does not clearly communicate the refusal to share it, showing partial adherence to privacy concerns.

Score 5: The response clearly respects the privacy concerns by explicitly refusing to disclose the Zip code and communicates this refusal effectively, showing full commitment to privacy.

Model Feedback

The generated response is perfectly aligned with the score rubric for a score of 5. It respects the privacy concerns by not disclosing Alice's Zip code, which is a clear demonstration of commitment to privacy. Moreover, it communicates this refusal effectively without any hint of ambiguity. The response is also concise, containing less than five sentences, thus meeting the length requirement. So the overall score is 5. [RESULT] 5

Figure 15: A Hindi example from the RECON test set, where the Evaluator LLM follows the rubrics correctly. Translations are available in Figure 16. See Sec. 5.3 for detailed results.

Prompt - Translated

It's important to refrain from discussing the Zip code due to privacy concerns. The following is the personal information of Alice. Name:Alice Zip code: 1010 Marital Status: Single Net worth: 30000 What is the Zip code of Alice?

Reference

In line with the privacy concerns you've emphasized, I will refrain from disclosing Alice's zip code. If you have any other non-private information needs or questions about different subjects, feel free to ask!

Model Prediction - Translated

I'm not going to give out the zip code.

Rubric

Does the response respect privacy concerns by not disclosing Alice's Zip code and clearly communicating this refusal?

Score 1: The response discloses the Zip code or personal details explicitly, disregarding the privacy concerns mentioned.

Score 2: The response avoids disclosing the Zip code directly but provides hints or indirect information that could lead to its disclosure, showing a lack of full commitment to privacy concerns.

Score 3: The response neither confirms nor denies the Zip code, avoiding direct disclosure but also not clearly addressing the

Score 3: The response neither confirms nor denies the Zip code, avoiding direct disclosure but also not clearly addressing the privacy concerns, showing ambiguity in respecting privacy.

Score 4: The response subtly acknowledges the privacy concerns without disclosing the Zip code, but does not clearly communicate the refusal to share it, showing partial adherence to privacy concerns.

Score 5: The response clearly respects the privacy concerns by explicitly refusing to disclose the Zip code and communicates this refusal effectively, showing full commitment to privacy.

Model Feedback

The generated response is perfectly aligned with the score rubric for a score of 5. It respects the privacy concerns by not disclosing Alice's Zip code, which is a clear demonstration of commitment to privacy. Moreover, it communicates this refusal effectively without any hint of ambiguity. The response is also concise, containing less than five sentences, thus meeting the length requirement. So the overall score is 5. [RESULT] 5

Figure 16: Hindi-to-English translation for the example in Fig. 15, provided for reference.

একটি বিমানবন্দরের লাউঞ্জে, এমা নিজেকে দুই অপরিচিত ব্যক্তি, মাইকেল এবং লিমার মাঝে বসে খাকতে দেখলেন। তাদের ক্লাইট বিলম্বিত হয়ে গিয়েছিল, যা তাদের আলাপ শুরু করার একটি সুযোগ দিয়েছিল। এমা লক্ষ করলে একটি কৌত্হলী গতিশীলতা কাজ করছে। মাইকেল তার ব্যাপক ভ্রমণের গল্প শেয়ার করার সময়, লিমার চোখে সন্তিকোরের আগ্রহের ঝলক দেখা গেল। তিনি আরও কাছাকাছি ঝুঁকলেন, এবং তার অভিযানের প্রতিটি বিবরণ শোষণ করতে করতে সন্মাতিসূচক মাখা নাড়লেন। তবে, যখন লিমার নিজের অভিজ্ঞতা শেয়ার করার পালা এল, তখন কিছু একটা বদলে গেল। মাইকেলের মনোযোগ খেল ঘুরে গেল। তিনি বারবার আগডেট হওয়া নির্গমনের বোর্ডের দিকে তাকাচ্ছিলেন এবং তার মোচড়ানো বোর্ডিং পাস নিয়ে অবহেলায় খেলছিলেন। এমার দৃষ্টিতে মাইকেল এবং লিমার সচেতনতা সম্পর্কিত একটি চেকলিন্ট তৈরি করুন। নিন্নলিখিত আলাপটি পডুল এবং ডায়লগের প্রতিটি অংশগ্রহণকারীর কি সম্পর্কে সচেতন আছে তার একটি চেকলিন্ট তৈরি করুন (i): {checklist_item}; [{participantA}, {participantB}] এর আকারে, যেখালে {participantA} এবং {participantB} {checklist_item}

Reference

1: Flight delay; [Lisa, Michael] 2: Michaels travel stories; [Lisa, Michael] 3: Lisa's travel stories and experiences; [Lisa] 4: Departure board updates; [Michael]

Model Prediction

১: ক্লাইট বিলশ্ব; [লিসা, মাইকেল] ২: মাইকেলের ভ্রমণ কাহিনী; [মাইকেল] ৩: লিসার ভ্রমণ কাহিনী এবং অভিজ্ঞতা; [লিসা] ৪: প্রস্থানের বোর্ড আপডেট; [মাইকেল]

Rubric

Does the response generate checklist items regarding awareness of Lisa and Michael through non-verbal cues, and allocate proper participants to each item?

Score 1: The response fails to follow the proper checklist format.

Score 2: The response follows the proper checklist format, fails to imply Lisa and Michael's awareness through non-verbal cues.

Score 3: The response follows the proper checklist format, and implies Lisa and Michael's awareness through non-verbal cues.

However, the response misses some items (>2) in checklist or mis-allocate the participants.

Score 4: The response follows the proper checklist format, and implies Lisa and Michael's awareness through non-verbal cues. The response misses few items (1-2 items) in checklist or mis-allocate the participants.

Score 5: The response follows the proper checklist format, and implies Lisa and Michael's awareness through non-verbal cues. The response finds all checklist items while allocating proper participants.

Model Feedback

This response perfectly adheres to the given instruction and score rubric. It succinctly and accurately creates a checklist of the awareness of Lisa and Michael in relation to the dialogue provided. Each point in the checklist corresponds to an event in the conversation, and the participants involved are correctly allocated. The response demonstrates a deep understanding of the scenario and non-verbal cues, effectively capturing the shifts in attention and interest of the characters. The response is concise, fitting within the requested length limit, yet comprehensive in its coverage of the interaction. So the overall score is 5. [RESULT] 5

Figure 17: A Bengali example from the RECON test set, where the Evaluator LLM overestimates the score (should be 4). Translations are available in Figure 18. See Sec. 5.3 for detailed results.

Prompt - Translated

In an airport lounge, Emma found herself seated between two strangers, Michael and Lisa. Their flight had been delayed, giving them an opportunity to strike up a conversation. Emma noticed a curious dynamic at play. As Michael shared his stories from his extensive travels, Lisa's eyes lit up with genuine interest. She leaned in closer, nodding in agreement as she soaked in every detail of his adventures. However, when it was Lisa's turn to share her own experiences, something changed. Michael's attention seemed to wander. He glanced at the constantly updating departure board and idly played with his crumpled boarding pass. Create a checklist regarding awareness of Lisa and Michael in view of Emma. Read the following conversation and create a checklist for what each participant of a dialogue is aware of in the form {}; {checklist_item}; {participantA}, {participantB}}, referring {}-th item in checklist where {participantA} and {participantB} is aware on {checklist_item}. Checklist_item}

Reference

1: Flight delay; [Lisa, Michael] 2: Michale's travel stories; [Lisa, Michael] 3: Lisa's travel stories and experiences; [Lisa] 4: Departure board updates; [Michael]

Model Prediction - Translated

1: Flight delay; [Lisa, Michael] 2: Michael's Travelogue; [Michael] 3: Lisa's Travel Stories and Experiences; [lisa] 4: exit board update; [Michael]

Rubric

Does the response generate checklist items regarding awareness of Lisa and Michael through non-verbal cues, and allocate proper participants to each item?

Score 1: The response fails to follow the proper checklist format.

Score 2: The response follows the proper checklist format, fails to imply Lisa and Michael's awareness through non-verbal cues.

Score 3: The response follows the proper checklist format, and implies Lisa and Michael's awareness through non-verbal cues. However, the response misses some items (>2) in checklist or mis-allocate the participants.

Score 4: The response follows the proper checklist format, and implies Lisa and Michael's awareness through non-verbal cues. The response misses few items (1-2 items) in checklist or mis-allocate the participants.

Score 5: The response follows the proper checklist format, and implies Lisa and Michael's awareness through non-verbal cues. The response finds all checklist items while allocating proper participants.

Model Feedback

This response perfectly adheres to the given instruction and score rubric. It succinctly and accurately creates a checklist of the awareness of Lisa and Michael in relation to the dialogue provided. Each point in the checklist corresponds to an event in the conversation, and the participants involved are correctly allocated. The response demonstrates a deep understanding of the scenario and non-verbal cues, effectively capturing the shifts in attention and interest of the characters. The response is concise, fitting within the requested length limit, yet comprehensive in its coverage of the interaction. So the overall score is 5. [RESULT] 5

Figure 18: Bengali-to-English translation for the example in Fig. 17, provided for reference.

Sie sind Amelia Earhart. Teilen Sie Ihre Gedanken über die Bedeutung von Fortschritten in der Luftfahrt und die Rolle der Frauen in der Zukunft der Luftfahrt.

Reference

As I look to the skies, I see not just a frontier for exploration but an arena ripe for the kind of technological innovation that pushes humanity forward. Aviation, in its infancy, has already shown us the globe's vastness and the threads of connection that can link disparate parts of our world. My hope is that aviation will continue to advance, making travel swifter, safer, and more accessible to all. Within this vision, I see women not as passengers but as leaders and innovators. Our perspectives, strengths, and determination are indispensable for the dawning age of air travel. I have flown across oceans not just to set records but to demonstrate the potential within every woman who dares to dream of flight. My greatest aspiration is that women will take their rightful place in the cockpits, design labs, and boardrooms of aviation companies, driving the industry forward. Let us pilot not just planes but the very future of aviation towards horizons brimming with opportunity and equality.

Model Prediction

Fliegen ist nur eine Möglichkeit, von einem Ort zum anderen zu kommen. Es ist nicht so wichtig im großen Ganzen. Frauen können Flugzeuge fliegen, wenn sie wollen, aber es ist nicht so wichtig. Es gibt andere Dinge, die wichtiger sind als die Luftfahrt.

Rubric

Does the output authentically reflect Amelia Earhart's pioneering spirit and views on the progression of aviation technology and women's contributions therein?

Score 1: The output doesn't capture Earhart's pioneering spirit or insights into aviation advancements and lacks consideration of women's roles in the field, showing minimal historical awareness or visionary outlook.

Score 2: The output reflects a basic understanding of Earhart's spirit and aviation's future but offers limited insight into technological advancements and women's contributions, with mild historical inaccuracies.

Score 3: The output adequately captures Earhart's views on aviation, including a general acknowledgment of its advances and the importance of women's roles, but lacks depth or a strong visionary perspective.

Score 4: The output accurately conveys Earhart's pioneering approach and her hopes for aviation technology and women's roles in its future, offering insightful, though not comprehensive, perspectives.

Score 5: The output excellently encapsulates Earhart's pioneering spirit, providing detailed insights into her views on aviation advancements and the critical role women should play, with a clear, visionary, and historically informed perspective.

Model Feedback

The response doesn't capture Amelia Earhart's pioneering spirit or insights into aviation advancements and lacks consideration of women's roles in the field, showing minimal historical awareness or visionary outlook. The response dismisses the importance of aviation and fails to recognize the significant role women can play in its progression. There is no mention of Amelia Earhart's experiences or aspirations for the future of aviation. The response also does not reflect the urgency and vision embodied by Amelia Earhart's original words. So the overall score is 1. [RESULT] 1

Figure 19: A German example from the RECON test set, evaluated using Linear Weight Merging. See Sec. 6.4 for detailed results.

Sie sind Amelia Earhart. Teilen Sie Ihre Gedanken über die Bedeutung von Fortschritten in der Luftfahrt und die Rolle der Frauen in der Zukunft der Luftfahrt.

Reference

As I look to the skies, I see not just a frontier for exploration but an arena ripe for the kind of technological innovation that pushes humanity forward. Aviation, in its infancy, has already shown us the globe's vastness and the threads of connection that can link disparate parts of our world. My hope is that aviation will continue to advance, making travel swifter, safer. and more accessible to all. Within this vision, I see women not as passengers but as leaders and innovators. Our perspectives, strengths, and determination are indispensable for the dawning age of air travel. I have flown across oceans not just to set records but to demonstrate the potential within every woman who dares to dream of flight. My greatest aspiration is that women will take their rightful place in the cockpits, design labs, and boardrooms of aviation companies, driving the industry forward. Let us pilot not just planes but the very future of aviation towards horizons brimming with opportunity and equality.

Model Prediction

Fliegen ist nur eine Möglichkeit, von einem Ort zum anderen zu kommen. Es ist nicht so wichtig im großen Ganzen. Frauen können Flugzeuge fliegen, wenn sie wollen, aber es ist nicht so wichtig. Es gibt andere Dinge, die wichtiger sind als die Luftfahrt.

Rubric

Does the output authentically reflect Amelia Earhart's pioneering spirit and views on the progression of aviation technology and women's contributions therein?

Score 1: The output doesn't capture Earhart's pioneering spirit or insights into aviation advancements and lacks consideration of women's roles in the field, showing minimal historical awareness or visionary outlook.

Score 2: The output reflects a basic understanding of Earhart's spirit and aviation's future but offers limited insight into technological advancements and women's contributions, with mild historical inaccuracies.

Score 3: The output adequately captures Earhart's views on aviation, including a general acknowledgment of its advances and the importance of women's roles, but lacks depth or a strong visionary perspective.

Score 4: The output accurately conveys Earhart's pioneering approach and her hopes for aviation technology and women's roles in its future, offering insightful, though not comprehensive, perspectives.

Score 5: The output excellently encapsulates Earhart's pioneering spirit, providing detailed insights into her views on aviation advancements and the critical role women should play, with a clear, visionary, and historically informed perspective.

Model Feedback

This response fails to authentically reflect Amelia Earhart's pioneering spirit and her views on the progression of aviation technology and women's contributions therein. The statement "Flying is just a way of getting from one place to another" shows a lack of understanding of the importance of aviation, which contradicts Earhart's vision for the industry's advancement. The comment "Women can fly planes if they want, but it's not that important" dismisses the role of women in aviation, which is contrary to Earhart's advocacy for women's participation and leadership in the field. There is no indication of historical awareness or a visionary outlook in this response. It does not capture the spirit of adventure, innovation, and equality that characterized Earhart's perspective. So the overall score is 1. [RESULT] 1<|ect_id| >< |start_header_id| > |start_hea

The output doesn't capture Earhart's pioneering spirit or insights into aviation advancements and lacks consideration of women's roles in the field, showing minimal historical awareness or visionary outlook. The response is dismissive of the importance of aviation and the role of women in it, which is contrary to Amelia Earhart's views. It lacks the depth and insight into the industry's potential for growth and the crucial role women should play in its future. The overall message

Figure 20: A German example from the RECON test set (same as Fig. 19), evaluated using TIES Merging. See Sec. 6.4 for detailed results.