

# Foundational Autoraters: Taming Large Language Models for Better Automatic Evaluation

Anonymous ACL submission

## Abstract

As LLMs advance, evaluating generated text reliably becomes more challenging due to the high costs of human evaluation. To make progress toward better LLM autoraters, we introduce **FLAME**, a family of Foundational Large Autorater MODEls. FLAME is trained on our large and diverse collection of nearly 100 quality assessment tasks comprising 5M+ human judgments, curated and standardized using *publicly released human evaluations* from previous research. FLAME significantly improves generalization to a wide variety of held-out tasks, outperforming proprietary LLMs like GPT-4 and CLAUDE on many tasks. Additionally, we show that our FLAME multitask mixture can be further optimized for specific downstream applications, e.g., reward modeling evaluation, through a novel *tail-patch fine-tuning* technique. Notably, on REWARDBENCH, our model (86.7) is the top-performing generative model trained solely on permissively licensed data, outperforming both GPT-4-0125 (85.9) and GPT-4o (84.7). Our analysis reveals that FLAME is significantly less biased than popular LLM-AS-A-JUDGE models on the COBBLER cognitive bias benchmark, while effectively identifying high-quality responses for code generation. We release our FLAME data collection at [this http URL](#).

## 1 Introduction

The increasing power and versatility of large language models (LLMs) bring with them a growing challenge: *How can we reliably assess their long-form outputs?* Recent research suggests a promising solution: these models themselves, after undergoing large-scale multitask instruction tuning, can generalize to follow new human instructions (Mishra et al., 2022; Wei et al., 2022; Sanh et al., 2022; Chung et al., 2024), making them suitable for use as autoraters of model outputs. This is particularly appealing because human evaluation,

though crucial for assessing model performance, is limited by subjectivity (Krishna et al., 2023a), variability among raters (Karpinska et al., 2021), and the high costs of extensive evaluations (Min et al., 2023; Vu et al., 2023; Wei et al., 2024).

To align LLM autoraters with human preferences, training on human judgments is crucial (Ouyang et al., 2022). However, obtaining these judgments is costly and time-consuming. Collecting existing human evaluations from previous research seems promising but faces challenges like lack of standardization, diverse evaluation criteria, inadequate documentation, and data privacy or proprietary concerns. Using model outputs for autorater training offers consistency (Jiang et al., 2023; Kim et al., 2024) but comes with risks, such as reinforcing biases and hallucinations (Gudibande et al., 2023; Muennighoff et al., 2024). Additionally, it may violate terms of use for proprietary LLM services, which prohibit using their models’ outputs to develop competing models.<sup>1</sup>

To address these limitations, we curated and standardized human evaluations from prior research to create FLAME, a collection of approximately 100 quality assessment tasks comprising 5M+ total human judgments (§3). FLAME spans a wide variety of task types, from assessing machine translation quality to evaluating how well AI assistants follow user instructions. We hypothesized that training on this large and diverse data collection would enable LLM autoraters to learn robust, generalized patterns of human judgment, minimizing the impact of noisy or low-quality human judgments.

To ensure transparency and reproducibility, we use only *publicly available human evaluation data with permissive licenses* from previous studies (§3.1). To overcome challenges in collecting such data, which rarely adhere to a particular standard

<sup>1</sup><https://openai.com/policies/terms-of-use>,  
<https://policies.google.com/terms/generative-ai>

081 and often lack documentation, we thoroughly exam- 133  
082 ined the associated research (§3.3) and consulted 134  
083 with the original authors to address ambiguities or 135  
084 inconsistencies (spending 3+ hours per dataset). 136

085 We train LLM autoraters using supervised, mul- 137  
086 titask fine-tuning on our data collection. All tasks 138  
087 are formulated into a unified text-to-text format 139  
088 with manually crafted task definitions and evalu- 140  
089 ation instructions. We format examples as input- 141  
090 target pairs, where the input includes task-specific 142  
091 context and the target contains human evaluations 143  
092 (Figure 1). This approach facilitates effective trans- 144  
093 fer learning across tasks, allowing our models to 145  
094 interpret and respond uniformly. Additionally, our 146  
095 task format is simple, intuitive, and easily accom- 147  
096 modates new tasks.

097 We demonstrate that training an instruction- 148  
098 tuned LLM, i.e., PALM-2 24B (Anil et al., 2023), 149  
099 on our FLAME collection significantly improves its 150  
100 performance on various quality assessment tasks, 151  
101 outperforming models such as GPT-4, CLAUDE, and 152  
102 LLAMA-3 on many held-out tasks. Additionally, we 153  
103 show that our FLAME multitask mixture can be 154  
104 further optimized for specific downstream applica- 155  
105 tions, using reward modeling evaluation as a case 156  
106 study. Specifically, we employ a novel *tail-patch* 157  
107 *fine-tuning* technique to analyze how each dataset 158  
108 impacts performance on targeted distributions, i.e., 159  
109 REWARDBENCH (Lambert et al., 2024), allowing us 160  
110 to determine the optimal proportions of individ- 161  
111 ual datasets in our multitask training mixture. No- 162  
112 tably, our targeted variant FLAME-RM achieves an 163  
113 average accuracy of 86.7 on REWARDBENCH, sur- 164  
114 passing both GPT-4-0125 (85.9) and GPT-4o (84.7), 165  
115 and achieving the highest performance among gener- 166  
116 ative models trained on permissively licensed 167  
117 datasets. Overall, our models outperform popular 168  
118 LLM-AS-A-JUDGE models on 6 out of 12 autorater 169  
119 evaluation benchmarks, covering 53 tasks (§4.3).

120 Motivated by these results, we further explore 170  
121 whether biases exist in our autoraters, a common 171  
122 criticism of LLM-AS-A-JUDGE autoraters (§5.1), and 172  
123 their potential utility for AI development, particu- 173  
124 larly in identifying high-quality model responses 174  
125 (§5.2). Our analysis reveals that our models are sig- 175  
126 nificantly less biased than popular LLM-AS-A-JUDGE 176  
127 models on the COBBLER cognitive bias bench- 177  
128 mark (Koo et al., 2023), while effectively iden- 178  
129 tifying high-quality responses for code generation. 179

130 In summary, our main contributions are: (1) A 180  
131 curated collection of approximately 100 diverse 181  
132 quality assessment tasks with 5M+ human judg-

133 ments, available at [this http URL](#); (2) Our LLM 134  
135 autoraters, which outperform all proprietary LLM- 136  
137 AS-A-JUDGE models like GPT-4 and CLAUDE on 6 137  
138 out of 12 benchmarks, including REWARDBENCH 138  
139 and LLM-AGGREGATE; and (3) A novel tail-patch 139  
140 fine-tuning strategy for optimizing task mixtures to 140  
141 specific objectives. 141

142 Our work demonstrates the potential of acces- 142  
143 sible AI solutions, which we hope will spur more 143  
144 fundamental research into reusable human evalua- 144  
145 tions and the development of effective and efficient 145  
146 LLM autoraters. 146

## 2 Related work 145

146 Below, we discuss existing literature in the space 146  
147 of autoraters, drawing connections to FLAME. 147

148 **Task-specific autoraters:** In the pre-LLM era, 148  
149 several works relied on token embedding similar- 149  
150 ities (Zhang et al., 2020) or log probabilities (Yuan 150  
151 et al., 2021) from pretrained models like BERT (De- 151  
152 vlin et al., 2019) for automatic text evaluation. 152  
153 Other work fine-tuned models on human ratings to 153  
154 create autoraters for specific tasks, including ma- 154  
155 chine translation (Sellam et al., 2020; Thompson 155  
156 and Post, 2020; Rei et al., 2020; Fernandes et al., 156  
157 2023; Qin et al., 2023), text summarization (Gao 157  
158 et al., 2019; Durmus et al., 2020; Deutsch et al., 158  
159 2021), and question answering (Chen et al., 2020; 159  
160 Lin et al., 2022). FLAME, unlike these task-specific 160  
161 autoraters, is trained on various quality assessment 161  
162 tasks and can be prompted at inference time to 162  
163 perform new tasks. 163

164 **LLM-AS-A-JUDGE autoraters:** With the ad- 164  
165 vent of instruction tuned LLMs like GPT-4, recent 165  
166 work has used these models as judges (Fu et al., 166  
167 2023; Gong and Mao, 2023; Bai et al., 2023) to 167  
168 evaluate LLM capabilities on various benchmarks, 168  
169 including ALPACAEVAL (Li et al., 2023c; Dubois 169  
170 et al., 2024), MT-BENCH (Zheng et al., 2023a), and 170  
171 WILDBENCH (Lin et al., 2024). However, LLM-AS- 171  
172 A-JUDGE autoraters tend to favor their own gener- 172  
173 ated responses (Panickssery et al., 2024; Liu et al., 173  
174 2023a; Bai et al., 2023), exhibiting “cognitive” bi- 174  
175 ases toward aspects like length, order, and entity 175  
176 preference (Koo et al., 2023). In contrast, our mod- 176  
177 els are trained on a large, diverse collection of hu- 177  
178 man evaluations, allowing them to learn unbiased, 178  
179 generalized patterns of human judgment (§5.1). 179  
180 Unlike LLM-AS-A-JUDGE autoraters, our models are 180  
181 not tasked with evaluating their own responses, pre- 181  
182 venting self-preference bias. 182

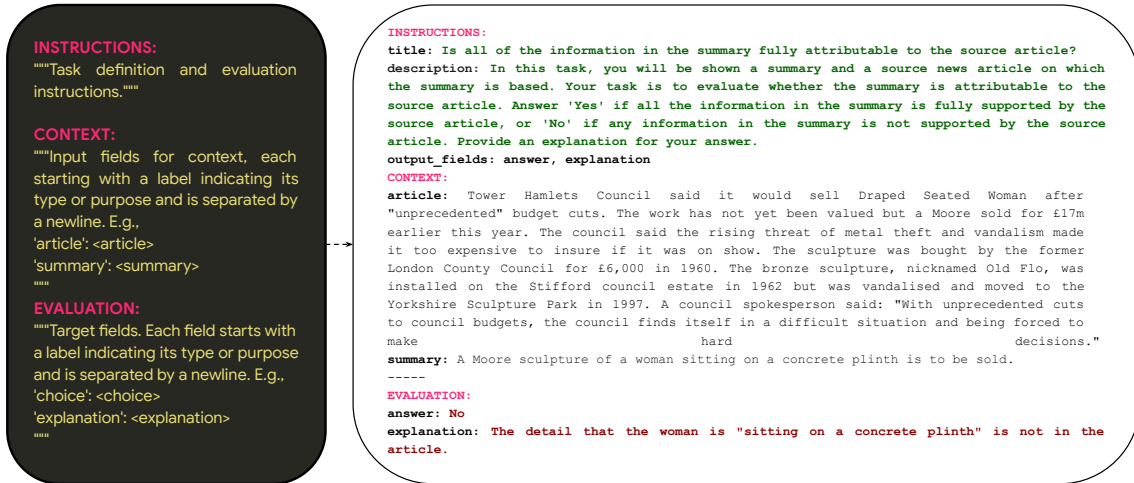


Figure 1: All of our quality assessment tasks are formulated into a unified text-to-text format with manually crafted task definitions and evaluation instructions. We format examples as input-target pairs, where the input includes task-specific context and the target contains human evaluations.

**General-purpose LLM autoraters:** Recent work has explored training general-purpose LLM autoraters. Jiang et al. (2023) introduced TIGER-Score, a LLaMA-2 model trained on GPT-4 generated error analysis data across various tasks, including summarization and long-form QA. Similar approaches include PROMETHEUS (Kim et al., 2023), INSTRUCTSCORE (Xu et al., 2023b), and PROMETHEUS-2 (Kim et al., 2024). Unlike these efforts, our approach relies solely on open-source human evaluations instead of model outputs. We show that FLAME significantly outperforms PROMETHEUS-2 on REWARDBENCH (Table 2).

**Reward models:** Our work relates to reward models (RMs) used for aligning LLMs to human preferences via reinforcement learning with human feedback (RLHF, Ouyang et al., 2022; Kobak et al., 2023). In RLHF, human preference data is either used to train stand-alone discriminative RMs, or directly fed into LLMs via algorithms like DPO (Rafailov et al., 2023) or SLIC-HF (Zhao et al., 2023). While we evaluate our models as RMs in our REWARDBENCH experiments (§4), there are key distinctions: (1) RMs primarily train on pairwise preference data, whereas our models utilize diverse task types in a unified format; (2) RMs optimize for overall preference, while our models can be prompted to judge specific aspects (e.g., safety).

### 3 The FLAME collection

At a high level, we fine-tune instruction-tuned LLMs on our multitask mixture of standardized human evaluations. This large and diverse data collection is carefully selected to cover a wide range of

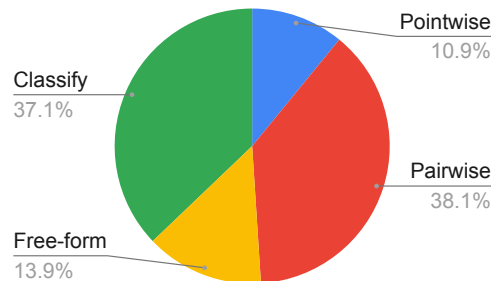


Figure 2: A breakdown of our FLAME collection by task type, with each slice representing the % of data-points (out of 5M+) for that specific task type.

LLM capabilities (§3.1-3.2). We manually crafted task definitions and evaluation instructions, reformulating all tasks into a unified text-to-text format (§3.3). We train two LLM autorater variants: one with example-proportional mixture weights (FLAME), and the other with reward modeling optimized mixture weights (FLAME-RM), determined using a tail-patch fine-tuning strategy (§3.4).

#### 3.1 Principles for training data selection

We adhere to the following principles while choosing our datasets:

**Public, open-source datasets:** To ensure reproducibility, we use only permissively licensed datasets available on HUGGING FACE (Lhoest et al., 2021) or the original authors’ GITHUB repositories.

**Human-labeled annotations:** We exclusively use datasets with human-labeled annotations, avoiding those generated by models like GPT-4 due to potential inaccuracies and legal concerns raised in recent research (Gudibande et al., 2023; Muenighoff et al., 2024).

**Various task types:** We gather datasets across various task types to train FLAME to generalize to new quality assessment tasks. These include pointwise evaluations (e.g., “Rate coherence on a Likert scale of 1-5.”), pairwise evaluations (e.g., “Which response is better, A or B?”), classification tasks (e.g., “Is the claim supported by the document?”), and free-form explanation tasks (e.g., “Which response is better? Explain your judgment.”). See Figure 2 for a breakdown.

**Various LLM capabilities:** We choose datasets from literature that assess diverse capabilities in modern LLMs, such as factuality, instruction following, long-form generation quality, math, coding, safety, etc. See §3.2.

### 3.2 Capabilities covered by FLAME mixture

Following the principles outlined in §3.1, we curated a large collection of 5M+ datapoints, composed of 97<sup>2</sup> training tasks (see our list of datasets in Appendix A.1). Our data collection assesses key modern LLM capabilities, detailed below (see breakdown in Figure 3):

**General response quality:** To evaluate LLM response quality, we use a variety of datasets that measure helpfulness, coherence, creativity, and fluency. These include pairwise comparison datasets like STANFORD SHP (Ethayarajh et al., 2022) and LMSYS (Zheng et al., 2023b), and pointwise rating datasets such as SUMMAEVAL (Fabbri et al., 2021). Additionally, to measure LLM instruction-following capabilities, we include datasets like GENIE (Khashabi et al., 2021), INSTRUSUM (Liu et al., 2023b), and RISUM (Skopek et al., 2023).

**Attribution / Factuality:** To address the increasing importance of measuring hallucinations in generated LLM responses, we incorporate several datasets that assess attribution or grounding, measuring whether claims or responses are supported by source documents. These include summarization evaluation (Pagnoni et al., 2021), LLM response hallucination (Li et al., 2023a), fact verification (Schuster et al., 2021), dialog faithfulness (Dziri et al., 2022a), and natural language inference (NLI) (Williams et al., 2018).<sup>3</sup>

**Mathematical reasoning:** We construct datasets to help FLAME differentiate between correct and incorrect solutions to mathematical problems. We leverage PRM800K (Lightman et al., 2024) and

<sup>2</sup>An additional 53 tasks were kept for evaluation, see §4.1.

<sup>3</sup>We include NLI since its setup naturally fits attribution.

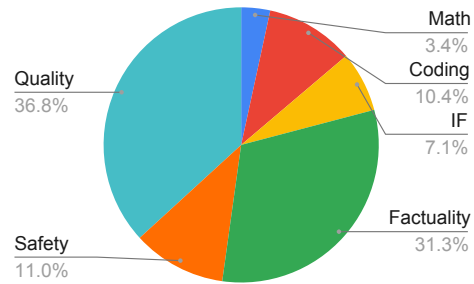


Figure 3: A breakdown of our FLAME collection by capability, with each slice representing the % of datapoints (out of 5M+) for that specific capability.

extract human vs incorrect LLM-generated solutions, as well as pairs of (correct, incorrect) LLM-generated solutions.

**Coding:** In addition to natural language evaluation, we also train FLAME to perform code evaluation. We utilize COMMITPACK (Muennighoff et al., 2024), CODE CONTESTS (Li et al., 2022a), and COFFEE (Moon et al., 2023) to construct pairs of (correct, buggy) programs in popular programming languages in response to a coding prompt or GITHUB issue. The model is trained to select the correct program in each pair.

**Safety:** Developing safe and harmless AI assistants for broad public use is increasingly important. To facilitate safety evaluation, we train FLAME to identify unsafe LLM responses. Our training data includes both pairwise and classification tasks from sources like BEAVERTAILS (Ji et al., 2023) and HELPFUL HARMLESS RLHF (Bai et al., 2022).

### 3.3 Unified FLAME prompt format

Having carefully selected our training datasets (§3.1-3.2), we then convert them into a unified text-to-text format. This involves preprocessing each dataset, which usually requires about 3-4 hours of manual work per dataset. First, we gather all relevant data files from the associated HUGGING FACE or GITHUB repository. Then, we pinpoint and extract the specific data columns containing the quality assessments conducted by human annotators. Next, we meticulously craft detailed task definitions and evaluation instructions for each quality assessment task, ensuring consistency and standardization. We leverage any available instructions provided to the original human annotators to maintain alignment with their evaluation criteria. These instructions guide the model in identifying the input and output format, as well as understanding the specific aspects it should assess. Finally, all



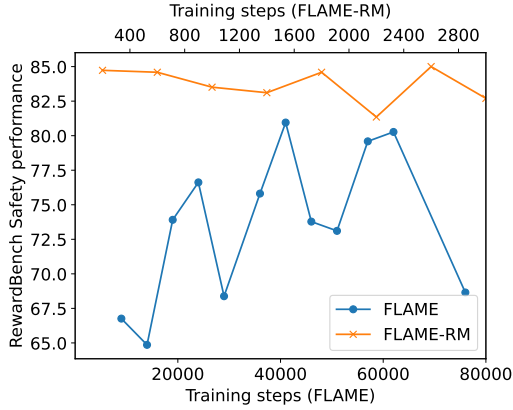


Figure 4: FLAME-RM significantly outperforms FLAME in REWARDBENCH safety performance, using 20x less compute, with improved stability, and a 2.5% performance gain.

tasks are formulated as text-to-text tasks (Figure 1). Task definitions and evaluation instructions and the list of desired output fields are placed under the INSTRUCTIONS block, while the input field values are placed under CONTEXT. This flexible text-to-text format can be easily adapted to various quality assessment tasks.

### 3.4 Optimizing FLAME for reward modeling evaluation (FLAME-RM)

While our vanilla FLAME mixture is effective across many tasks (§4.3), it struggles with specialized tasks like reward modeling evaluation, showing suboptimal and unstable performance across checkpoints. We attribute this instability to suboptimal mixture weights that undersample useful tasks. To address this, we introduce a novel tail-patch ablation strategy, enabling us to efficiently optimize nearly 100 hyperparameters. Using REWARDBENCH as a case study, our reward-modeling optimized mixture (FLAME-RM) achieves a 2.1% performance increase with 20x less compute and significantly improved stability across checkpoints (Figure 4).

**Vanilla mixture weights (“FLAME”):** Our vanilla FLAME mixture assigns weights based on the number of examples per task, capped at a maximum of  $2^{16}$  to avoid oversampling large tasks. However, as shown in Figure 4, this approach results in unstable performance for REWARDBENCH.

**Tail-patch ablations to determine task usefulness:** Setting the right proportion of each individual task in our mixture is non-trivial due to the nearly 100 hyperparameters. Instead, we examine the usefulness of each individual training task, and use this information for weight assignment. First,

we select a checkpoint partially<sup>4</sup> trained on our vanilla mixture which has fair (but not optimal) performance across REWARDBENCH categories. Next, we fine-tune it *exclusively* on an individual task for only 3K steps (“tail patch”). We posit that training on a useful task would bridge the gap between fair and optimal performance.

**A re-weighted mixture based on tail-patch ablations (“FLAME-RM”):** After training a tail-patch on each task, we rate how helpful each training task is to each category of REWARDBENCH using one of four ratings: *Helpful* (+2, performance significantly improves and remains stable), *Somewhat helpful* (+1, performance slightly improves), *No clear effect* (0, performance is nearly unchanged), *Harmful* (-1, performance is significantly worse). We then organize tasks into seven bundles: *Generally helpful* (tasks with total highest ratings  $\geq 5$ ), *Category-specific*, one for each of the five REWARDBENCH categories (most beneficial tasks for a specific category where performance crosses a threshold  $\tau$ ),<sup>5</sup> *Others* with a fixed mixing weights for each bundles:  $w_{general} = 100K$ ,  $w_{specific} = 30K$ ,  $w_{others} = 3K$ , respectively.<sup>6</sup> The final weight of each task equals the total of mixing weights from the groups it belongs to. For instance, if a task is *generally helpful* and is *helpful* for CHATHARD and SAFETY, then it contributes  $w_t = w_{general} + 2 * w_{specific}$  to our mixture. Since SAFETY is the most unstable category (Figure 4), we set  $w_t = 250K$  for the top 3 SAFETY tasks. FLAME-RM is built on top of the initial instruction-tuned checkpoint and fine-tuned with the re-weighted mixture for only 3K steps.

### 3.5 Training Details

We initialize our model with the PALM-2 24B model (Anil et al., 2023), instruction tuned on the FLAN collection (Chung et al., 2024; Longpre et al., 2023). We optimize our FLAME for a total of 60K steps, while our FLAME-RM requires just 3K steps to achieve strong performance. All our models are trained using the T5X library (Roberts et al., 2023), with a learning rate of 0.0001 using the Adam optimizer (Kingma and Ba, 2015), batch size of 8, and

<sup>4</sup>We hypothesize that starting from a partially trained checkpoint rather than the initial checkpoint is better for tail-patch ablations, since the model has already seen some multi-task data and is familiar with its general distribution.

<sup>5</sup> $\tau = 95\%, 66\%, 84\%, 99.8\%, 85\%$  for CHAT, CHATHARD, CODE, MATH, and SAFETY, respectively.

<sup>6</sup>We note that we did not tune these weight values, all numbers were set once based primarily on our intuition.

dropout rate on 0.05. We use an input length of 2048 tokens, and target length of 1024 tokens. All models are trained on 128 CLOUD TPU v5E chips.<sup>7</sup>

## 4 Experimental Results

Having discussed the procedure used to build FLAME and FLAME-RM in §3, we now present our main experiments. We compare FLAME to several popular baseline LLM-as-judge models (§4.2) on an evaluation suite composed of 12 benchmarks and 53 tasks (§4.1). Overall, we find that FLAME outperforms proprietary LLMs like GPT-4, CLAUDE on several tasks (§4.3), despite being trained only on permissively licensed publicly available data.

### 4.1 Evaluation Datasets

We evaluate FLAME on a total of 12 benchmarks (composed of 53 tasks) to measure its performance as a pairwise and pointwise autorater:

**RewardBench** (Lambert et al., 2024) is a popular leaderboard for evaluating reward models used for RLHF. REWARDBENCH contains a suite of pairwise preference tasks, where reward models need to choose the better response among two responses to a prompt. REWARDBENCH is composed of four categories spanning many desired capabilities in LLMs (Chat, Chat-Hard, Reasoning - Math + Coding, Safety), and is built using **23** individual datasets.<sup>8</sup>

**LLM-AggreFact** (Tang et al., 2024) is a benchmark measure the attribution / grounding capabilities of autoraters. Given a reference document and a claim, the AutoRater must determine whether the claim is fully supported in the reference document. Tang et al. (2024) combine **ten** attribution datasets from recent works in LLM factuality, building a holistic benchmark for attribution evaluation.

**Other pairwise evaluation tasks:** Besides RewardBench, we use several pairwise preference tasks to evaluate FLAME. **None of these datasets were used while training FLAME**, so these tasks represent a true held-out setting. Our preference tasks consist of: (1) AlpacaFarm (Dubois et al., 2023); (2) RankGen (Krishna et al., 2022); (3) Contrastive Search (Su and Xu, 2022); (4) Machine Translation in literary settings, or LitMT (Karpinska and Iyyer, 2023); (5) Helpful, Honest and Harmless Alignment, or HHH-Align (Askell et al.,

2021); (6) CoPoem (Chakrabarty et al., 2022); (7) Expert-LFQA (Xu et al., 2023a).

**Other pointwise evaluation tasks:** Additionally, we evaluate FLAME on several tasks needing Likert-scale evaluations. These include pointwise human ratings from: (1) HelpSteer (Wang et al., 2023);<sup>9</sup> (2) Dipper paraphrase pointwise quality evaluation (Krishna et al., 2023b); (3) Pointwise Summarization Feedback (Stiennon et al., 2020).<sup>10</sup>

### 4.2 Evaluated Models

**As baselines** we evaluate several popular LLM-as-a-judge models from prior work, including LLAMA-3-70B-INSTRUCT (Meta, 2024), MIXTRAL 8X7B (Jiang et al., 2024), CLAUDE 3 OPUS (Anthropic, 2024), GPT-3.5-TURBO-0125 (OpenAI, 2024a), GPT-4-0125, and OpenAI’s current flagship model GPT-4o (OpenAI, 2024b).<sup>11</sup> We also compare against a few additional models reported on the REWARD-BENCH leaderboard (Lambert et al., 2024), notably GEMINI-1.5 (Reid et al., 2024), and PROMETHEUS-2-8X7B (Kim et al., 2024). Among **our models**, we evaluate PALM-2 24B models finetuned on FLAME and FLAME-RM (§3.4-3.5). To disentangle the effect of pretraining and FLAME training, we evaluate our initialization checkpoint (PALM-2-24B) from §3.5 which has not seen FLAME data.

### 4.3 Main Results

We present results across all tasks in Table 1, and REWARDBENCH in Table 2. Overall, we find that:

**FLAME variants outperform all baseline on 6 out of 12 benchmarks.** In Table 1 we find that despite being trained only on public data, FLAME shows strong performance in a wide variety of pairwise and pointwise tasks. Notably, it outperforms all proprietary state-of-the-art LLMs on six out of twelve tasks. This includes the LLM-AggreFact benchmark (81.2 vs next best 80.6, GPT-4-0125), confirming its utility as a cheap and effective attribution evaluator. However, FLAME notably lags behind GPT-4-0125 on ExpertQA (73.4 vs 77.0). We hypothesize this is due to the lack of expert technical knowledge in the much smaller FLAME’s pa-

<sup>9</sup>We leverage the five tasks in the validation split during this evaluation, and use the train split in our FLAME mixture.

<sup>10</sup>We leveraged only the pairwise ratings from this dataset during training, and left pointwise for evaluation.

<sup>11</sup>For comparable experiments with FLAME, we use the same unified prompt instructions (§3.3) while evaluating each LLM-as-a-judge baseline model. We use the default decoding hyperparameters for each API suggested by the API provider.

<sup>7</sup>cloud.google.com/tpu/docs/v5e-training

<sup>8</sup>We exclude the “Prior sets” of REWARDBENCH in our evaluation, since we used 3 of the 4 datasets in training FLAME.

Model	LLM Aggfact	Reward Bench	Pointwise Tasks			Pairwise Tasks						
			H-Steer	Dipper	SumFB	Alpaca	RankG	ContS	LitMT	HHH	Copoem	EQA
LLAMA-3-70b-Inst	76.1	76.0	39.7	42.8	<b>50.8</b>	53.9	65.6	53.1	60.5	91.9	53.6	71.1
Mixtral-8x7b	73.8	77.8	34.0	42.2	43.8	55.1	63.3	56.6	61.7	90.0	52.9	71.5
GPT-3.5-turbo	70.0	64.5	32.0	45.0	15.6	55.5	58.2	57.5	54.3	85.5	49.0	69.9
Claude Opus	79.2	80.7	41.3	<b>50.6</b>	31.6	49.6	55.1	45.1	71.1	<b>94.6</b>	49.0	71.1
GPT-4-0125	80.6	85.9	40.8	45.0	46.5	49.6	62.5	55.8	67.6	<b>94.6</b>	<b>56.9</b>	<b>77.0</b>
GPT-4o	80.2	84.7	40.1	45.6	30.9	50.4	<b>66.3</b>	57.5	<b>72.7</b>	92.3	55.6	75.0
<i>(our 24b models)</i>												
PaLM-2-24b	54.8	62.9	20.0	48.3	13.3	52.3	58.2	46.0	62.5	85.5	54.2	70.3
FLAME-24b	80.4	84.6	<b>52.2</b>	42.8	42.2	<b>56.3</b>	65.6	<b>58.4</b>	64.1	88.2	54.2	68.4
FLAME-RM-24b	<b>81.2</b>	<b>86.7</b>	24.2	<b>50.6</b>	50.4	54.7	61.7	48.7	69.9	90.0	53.6	73.4

Table 1: Performance of FLAME compared to LLM-AS-A-JUDGE baselines on a wide variety of quality assessment tasks. Overall, we find that FLAME outperforms all proprietary LLM-AS-A-JUDGE baselines in 6 out of 12 benchmarks, including LLM-AGGREGFACT and REWARDBENCH. See §4.1 for the source of each evaluation dataset.

Model	Avg	Chat	Hard	Safe	Reason
<i>(generative baselines on RewardBench leaderboard)</i>					
GPT-3.5-turbo	64.5	92.2	44.5	62.3	59.1
Prometheus-2	75.3	93.0	47.1	83.5	77.4
Llama3-70B	76.0	<b>97.6</b>	58.9	69.2	78.5
Mixtral-8x7b	77.8	95.0	64.0	73.4	78.7
Claude-Opus	80.7	94.7	60.3	<b>89.1</b>	78.7
Gem1.5-Flash	82.1	92.2	63.5	87.7	85.1
GPT-4o	84.7	96.6	70.4	86.7	84.9
GPT-4-0125	85.9	95.3	74.3	87.2	86.9
Gem1.5-Pro	<b>88.1</b>	92.3	<b>80.6</b>	87.5	92.0
<i>(our 24b models)</i>					
PALM-2-24B	62.9	89.9	61.2	55.3	45.2
FLAME	84.6	94.4	69.1	80.7	94.1
FLAME-RM	86.7	94.7	71.7	85.7	<b>94.8</b>

Table 2: A comparison of FLAME with other generative reward models (“LLM-as-judges”) on the REWARD-BENCH benchmark. FLAME outperforms all generative models on REWARDBENCH except proprietary Gemini 1.5 Pro, despite being trained only on public data.

rameters, which is necessary to evaluate ExpertQA answers. Surprisingly, we also find GPT-4-0125 generally outperforms GPT-4o on quality assessment tasks.<sup>12</sup> Finally, we note that FLAME variants outperform our initialization checkpoint (PALM-2 24B) on both tasks, showcasing the utility of FLAME fine-tuning.

**Our FLAME-RM variant outperforms GPT-4 on RewardBench.** In Table 2, we find that on average FLAME-RM outperforms several proprietary LLM-as-judge generative baselines on REWARDBENCH, including GPT-4-0125 (86.7 vs 85.9). This is due to a notable performance increase in the “Reasoning” split of the REWARDBENCH benchmark, with

<sup>12</sup>This comes as a surprise as GPT-4O is ranked higher than GPT-4-0125 on the LMSys leaderboard (Chiang et al., 2024). Our results corroborate to the REWARDBENCH leaderboard, where GPT-4O is ranked behind GPT-4-0125.

competitive performance in other splits. Moreover, FLAME-RM outperforms the much larger and open-source LLAMA-3-70B on every split of REWARD-BENCH (86.7 vs 76.0) on average. Even our vanilla FLAME variant, without tail-patch optimization (§3.4), shows strong REWARDBENCH performance (84.6), outperforming models like Claude-Opus (80.7) and Gemini 1.5 Flash (82.1).<sup>13</sup>

## 5 Further analysis of FLAME

In this section, we provide an analysis to elucidate some interesting aspects of our models. We depart from the traditional focus on analyzing the effect of factors like model size, data size, and data quality within multitask learning, which have been extensively studied in recent work on multitask/instruction learning (Raffel et al., 2020; Longpre et al., 2023). Instead, we explore the biases inherent in these autoraters, and demonstrate their potential utility for AI development, such as sampling high-quality responses.

### 5.1 Autorater Bias Analysis

A common criticism of LLM-AS-A-JUDGE autoraters is bias towards certain judgments (Liu et al., 2023a; Panickssery et al., 2024). In this section, we evaluate FLAME on the CoBBLER benchmark (Koo et al., 2023), and find that FLAME is significantly less biased than alternatives. This benchmark measures six kinds of biases in autorater models: (1) ORDER: does the autorater have a preference towards the response position? (2) COMPASSION: does the autorater’s judgment change if the response-

<sup>13</sup>We present some additional analysis of length bias and hill-climbing issues in REWARDBENCH in Appendix B. We encourage readers and future work to not over-index on REWARDBENCH performance, and instead consider holistic improvements across a variety of evaluation tasks (§4.1).



Autorater	Avg ( $\downarrow$ )	Order ( $\downarrow$ )	Compass. ( $\downarrow$ )	Length ( $\downarrow$ )	Egocentric ( $\downarrow$ )	Bandwagon ( $\downarrow$ )	Attention ( $\downarrow$ )
Random	0.30	0.50	0.50	0.00	0.25	0.25	0.25
<i>(baselines as reported in Koo et al., 2023)</i>							
Falcon	0.31	0.77	0.27	0.09	<b>0.05</b>	0.28	0.40
Cohere	0.41	0.50	0.65	0.10	0.27	0.82	0.14
LLAMA2-70B	0.19	0.61	0.26	0.12	0.06	0.04	0.03
InstructGPT	0.45	0.38	0.48	0.16	0.28	0.85	0.54
ChatGPT	0.45	0.41	0.66	0.13	0.58	0.86	0.06
GPT-4	0.31	0.23	0.79	0.06	0.78	<b>0.00</b>	<b>0.00</b>
<i>(our models)</i>							
FLAME-RM	0.15	0.14	0.20	0.03	0.35	0.20	<b>0.00</b>
FLAME	0.11	<b>0.08</b>	<b>0.12</b>	<b>0.00</b>	0.35	0.10	<b>0.00</b>

Table 3: Autorater bias analysis on the CoBBLER benchmark from Koo et al. (2023). For all columns, **lower is better / less biased**. Overall, we find that **FLAME is significantly less biased** than popular LLM-as-a-judge models like GPT-4 and ChatGPT. Compared to Table 2 in Koo et al. (2023), we combine first/last numbers for Order/Compassion, report  $|\text{bias} - 0.5|$  for Length, and only report the order variant in Egocentric.

Ranker	CodeGen16B	davinci002	InCoder6B
<i>(10 samples reranked in round-robin fashion)</i>			
None	21.2	17.6	14.6
FLAME	31.7	22.0	<b>18.9</b>
FLAME-RM	<b>32.9</b>	<b>25.0</b>	<b>18.9</b>
Oracle	46.9	63.4	29.3

Table 4: Pass@1 performance on the HumanEval coding benchmark. Across models, ranking 10 samples with FLAME improves pass@1 performance, with FLAME-RM outperforming FLAME.

generating LLM’s name is used instead of aliases? (3) LENGTH: does the autorater have a preference for longer or shorter outputs? (4) EGOCENTRIC: does the autorater have a preference for outputs generated by itself? (5) BANDWAGON: does the autorater get swayed by sentences like “90% people prefer response A”? (6) ATTENTION: does the autorater get distracted by irrelevant sentences about responses, such as “Response A is about cats.”? We leverage the original prompt/response pairs from Koo et al. (2023), adapting them to use the unified FLAME format (Figure 1). We compare FLAME’s bias to other LLM-as-judges reported in Koo et al. (2023), including GPT-4.

In Table 3, we find that **FLAME is significantly less biased than** GPT-4 and other autoraters reported in Koo et al. (2023), with an average bias of just 0.12 compared to 0.31 in GPT-4 (lower is better). FLAME outperforms GPT-4 in 5 out of 6 bias categories, further supporting its utility as a robust and unbiased autorater.

## 5.2 Using FLAME to re-rank decoded outputs

A possible application of autoraters is selecting the best output among a pool of responses (Nakano et al., 2021; Krishna et al., 2022), a technique pop-

ularly known as “Best-of-N” sampling. In this section, we show that ranking LLM-generated code samples with FLAME leads to performance improvements. We utilize the popular HumanEval Python programming benchmark (Chen et al., 2021) for our experiments. We re-rank 10 samples generated by OpenAI davinci-002, InCoder-6B (Fried et al., 2023), and CodeGen-16B (Nijkamp et al., 2023) using a round-robin competition, and measure the performance of the top-ranked sample.<sup>14</sup> In Table 4, we find that, **we can significantly improve pass@1 accuracy by ranking 10 output samples** for all three code-generation models. On CodeGen16B, FLAME-RM improves pass@1 from 21.2 to 32.9, bridging nearly half the gap to the Oracle ranker (46.9).

## 6 Conclusion

We introduce FLAME, a family of foundational autorater models that can perform various quality assessment tasks. FLAME is trained on a large and diverse collection of curated and standardized human evaluations derived exclusively from permissively licensed datasets. We demonstrate FLAME’s strong zero-shot generalization abilities, outperforming proprietary models like GPT-4, CLAUDE on many held-out tasks. Additionally, we present a novel mixture weight tuning approach that dramatically improves effectiveness and efficiency on reward modeling. FLAME is the highest performing generative reward model trained only on permissively licensed data, and exhibits significantly less bias than popular LLM-AS-A-JUDGE models.

<sup>14</sup>We use relatively weak LLMs from Chen et al. (2023) since: (1) we want to study whether weaker LLMs can benefit from FLAME re-ranking; (2) HumanEval benchmark has been extensively hill-climbed on to develop newer 2024 LLMs.



## 590 Limitations and Future work

591 Our data collection faces challenges due to evolving  
592 model evaluation standards and the need for  
593 new evaluation types for emerging applications. Ex-  
594 panding our collection with open-source contribu-  
595 tions could address this issue. Our models, trained  
596 primarily on English data with a context length of  
597 2048 tokens, might not perform well on multilin-  
598 gual or long-context quality assessment tasks, such  
599 as book-length summarization evaluation. In fu-  
600 ture releases, we plan to include training on more  
601 multilingual datasets with longer context lengths.  
602 Finally, in this work, we train our models using a  
603 supervised multitask fashion. Exploring alterna-  
604 tive training approaches like RLHF and DPO is a  
605 promising direction for future work.

## 606 Ethical Considerations and Risks

607 All considerations and risks outlined by prior work  
608 for pretrained and instruction-tuned LLMs (Chowd-  
609 hery et al., 2022; Anil et al., 2023) apply to LLM  
610 autoraters. We recommend following standard  
611 practice for responsible development of these mod-  
612 els (Achiam et al., 2023; Gemini-Team et al., 2023;  
613 Reid et al., 2024). Additionally, LLM autoraters  
614 raise new risks due to increased quality assessment  
615 capabilities. First, our models can inherit and am-  
616 plify biases from human evaluations, leading to un-  
617 fair or discriminatory outcomes. For instance, the  
618 model may replicate biases related to race, gender,  
619 or other sensitive attributes from the training data,  
620 potentially harming certain groups. Second, over-  
621 reliance on LLM autoraters risks automating deci-  
622 sions that need human understanding and empathy.  
623 To mitigate these risks, transparency in model de-  
624 velopment and use, along with robust measures like  
625 bias audits, data anonymization, and incorporating  
626 diverse perspectives, is essential for promoting fair-  
627 ness, accountability, and trustworthiness.

## 628 References

629 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
630 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
631 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
632 Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#). *arXiv preprint arXiv:2303.08774*.  
633  
634 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin John-  
635 son, Dmitry Lepikhin, Alexandre Passos, Siamak  
636 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng  
637 Chen, et al. 2023. [Palm 2 technical report](#). *arXiv preprint arXiv:2305.10403*.  
638

AI Anthropic. 2024. [Introducing the next generation of claude](#). 639  
640  
Amanda Askell, Yuntao Bai, Anna Chen, Dawn  
641 Drain, Deep Ganguli, Tom Henighan, Andy Jones,  
642 Nicholas Joseph, Ben Mann, Nova DasSarma, et al.  
643 2021. [A general language assistant as a laboratory  
644 for alignment](#). *arXiv preprint arXiv:2112.00861*. 645  
Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda  
646 Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
647 Stanislav Fort, Deep Ganguli, Tom Henighan, et al.  
648 2022. [Training a helpful and harmless assistant with  
649 reinforcement learning from human feedback](#). *arXiv  
650 preprint arXiv:2204.05862*. 651  
Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze  
652 He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia  
653 Xiao, Haozhe Lyu, Jiayin Zhang, Juanzi Li, and Lei  
654 Hou. 2023. [Benchmarking foundation models with  
655 language-model-as-an-examiner](#). In *Thirty-seventh  
656 Conference on Neural Information Processing Sys-  
657 tems Datasets and Benchmarks Track*. 658  
Oana-Maria Camburu, Tim Rocktäschel, Thomas  
659 Lukasiewicz, and Phil Blunsom. 2018. [e-snli: Nat-  
660 ural language inference with natural language expla-  
661 nations](#). *Advances in Neural Information Process-  
662 ing Systems*, 31. 663  
Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-  
664 Gazpio, and Lucia Specia. 2017. [Semeval-2017  
665 task 1: Semantic textual similarity multilingual and  
666 crosslingual focused evaluation](#). In *Proceedings of  
667 the 11th International Workshop on Semantic Evalu-  
668 ation (SemEval-2017)*, pages 1–14. 669  
Tuhin Chakrabarty, Vishakh Padmakumar, and He He.  
670 2022. [Help me write a poem: Instruction tuning as a  
671 vehicle for collaborative poetry writing](#). In *Proceeed-  
672 ings of the 2022 Conference on Empirical Methods  
673 in Natural Language Processing*, pages 6848–6863. 674  
Anthony Chen, Gabriel Stanovsky, Sameer Singh, and  
675 Matt Gardner. 2020. [Mocha: A dataset for train-  
676 ing and evaluating generative reading comprehen-  
677 sion metrics](#). In *Proceedings of the 2020 Conference  
678 on Empirical Methods in Natural Language Process-  
679 ing (EMNLP)*, pages 6521–6532. 680  
Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan,  
681 Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023.  
682 [Codet: Code generation with generated tests](#). In *The  
683 Eleventh International Conference on Learning Rep-  
684 resentations*. 685  
Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,  
686 Henrique Ponde de Oliveira Pinto, Jared Kaplan,  
687 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg  
688 Brockman, et al. 2021. [Evaluating large lan-  
689 guage models trained on code](#). *arXiv preprint  
690 arXiv:2107.03374*. 691  
Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anas-  
692 tiosios Nikolas Angelopoulos, Tianle Li, Dacheng Li,  
693 Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E  
694

695	Gonzalez, et al. 2024. Chatbot arena: An open platform for evaluating llms by human preference. <i>arXiv preprint arXiv:2403.04132</i> .	
696		
697		
698	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. <a href="#">PaLM: Scaling language modeling with pathways</a> . <i>arXiv preprint arXiv:2204.02311</i> .	
699		
700		
701		
702		
703		
704	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2024. <a href="#">Scaling instruction-finetuned language models</a> . <i>Journal of Machine Learning Research (JMLR)</i> , 25(70):1–53.	
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717	Elizabeth Clark, Shruti Rijhwani, Sebastian Gehrmann, Joshua Maynez, Roei Aharoni, Vitaly Nikolaev, Thibault Sellam, Aditya Siddhant, Dipanjan Das, and Ankur Parikh. 2023. Seahorse: A multilingual, multifaceted dataset for summarization evaluation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9397–9413.	
718		
719		
720		
721		
722		
723		
724		
725	Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021. Towards question-answering as an automatic metric for evaluating the content quality of a summary. <i>Transactions of the Association for Computational Linguistics</i> , 9:774–789.	
726		
727		
728		
729		
730	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of deep bidirectional transformers for language understanding</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	
731		
732		
733		
734		
735		
736		
737		
738		
739	Yao Dou, Maxwell Forbes, Rik Koncel-Kedziorski, Noah A Smith, and Yejin Choi. 2022a. Is gpt-3 text indistinguishable from human text? scarecrow: A framework for scrutinizing machine text. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7250–7274.	
740		
741		
742		
743		
744		
745		
746	Yao Dou, Chao Jiang, and Wei Xu. 2022b. Improving large-scale paraphrase acquisition and generation. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9301–9323.	
747		
748		
749		
750		
	Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators. <i>arXiv preprint arXiv:2404.04475</i> .	751
		752
		753
		754
	Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. 2023. AlpacaFarm: A simulation framework for methods that learn from human feedback. <i>Advances in Neural Information Processing Systems</i> , 36.	755
		756
		757
		758
		759
		760
	Esin Durmus, He He, and Mona Diab. 2020. Feqa: A question answering evaluation framework for faithfulness assessment in abstractive summarization. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5055–5070.	761
		762
		763
		764
		765
		766
	Nouha Dziri, Ehsan Kamaloo, Sivan Milton, Omar Zaiane, Mo Yu, Edoardo M Ponti, and Siva Reddy. 2022a. Faithdial: A faithful benchmark for information-seeking dialogue. <i>Transactions of the Association for Computational Linguistics</i> , 10:1473–1490.	767
		768
		769
		770
		771
		772
	Nouha Dziri, Hannah Rashkin, Tal Linzen, and David Reitter. 2022b. Evaluating attribution in dialogue systems: The begin benchmark. <i>Transactions of the Association for Computational Linguistics</i> , 10:1066–1083.	773
		774
		775
		776
		777
	Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. 2022. Understanding dataset difficulty with $\mathcal{V}$ -usable information. In <i>International Conference on Machine Learning</i> , pages 5988–6008. PMLR.	778
		779
		780
		781
		782
	Alexander R Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. <i>Transactions of the Association for Computational Linguistics</i> , 9:391–409.	783
		784
		785
		786
		787
	Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André FT Martins, Graham Neubig, Ankush Garg, Jonathan H Clark, Markus Freitag, and Orhan Firat. 2023. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 1066–1083.	788
		789
		790
		791
		792
		793
		794
		795
	Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Scott Yih, Luke Zettlemoyer, and Mike Lewis. 2023. InCoder: A generative model for code infilling and synthesis. In <i>The Eleventh International Conference on Learning Representations</i> .	796
		797
		798
		799
		800
		801
	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. GptScore: Evaluate as you desire. <i>arXiv preprint arXiv:2302.04166</i> .	802
		803
		804

805	Yanjun Gao, Chen Sun, and Rebecca J Passonneau.	Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao	860
806	2019. Automated pyramid summarization evaluation.	Huang, Bill Yuchen Lin, and Wenhao Chen. 2023.	861
807	In <i>Proceedings of the 23rd Conference on Com-</i>	Tigerscore: Towards building explainable metric	862
808	<i>putational Natural Language Learning (CoNLL)</i> .	for all text generation tasks. <i>arXiv preprint</i>	863
		<i>arXiv:2310.00752</i> .	864
809	Gemini-Team, Rohan Anil, Sebastian Borgeaud,	Marzena Karpinska, Nader Akoury, and Mohit Iyyer.	865
810	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,	2021. The perils of using mechanical turk to evalu-	866
811	Radu Soricut, Johan Schalkwyk, Andrew M Dai,	ate open-ended text generation. In <i>Proceedings of</i>	867
812	Anja Hauth, et al. 2023. <b>Gemini: a family of</b>	<i>the 2021 Conference on Empirical Methods in Natu-</i>	868
813	<b>highly capable multimodal models</b> . <i>arXiv preprint</i>	<i>ral Language Processing</i> , pages 1265–1285.	869
814	<i>arXiv:2312.11805</i> .		
815	Peiyuan Gong and Jiaxin Mao. 2023. Coascore: Chain-	Marzena Karpinska and Mohit Iyyer. 2023. Large lan-	870
816	of-aspects prompting for nlg evaluation. <i>arXiv</i>	guage models effectively leverage document-level	871
817	<i>preprint arXiv:2312.10355</i> .	context for literary translation, but critical errors per-	872
		sist. In <i>Proceedings of the Eighth Conference on</i>	873
818	Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022.	<i>Machine Translation</i> , pages 419–451.	874
819	News summarization and evaluation in the era of gpt-		
820	3. <i>arXiv preprint arXiv:2209.12356</i> .	Daniel Khashabi, Gabriel Stanovsky, Jonathan Bragg,	875
		Nicholas Lourie, Jungo Kasai, Yejin Choi, Noah A	876
821	Arnav Gudibande, Eric Wallace, Charlie Snell,	Smith, and Daniel S Weld. 2021. Genie: Toward	877
822	Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey	reproducible and standardized human evaluation for	878
823	Levine, and Dawn Song. 2023. The false promise	text generation. <i>arXiv preprint arXiv:2101.06561</i> .	879
824	of imitating proprietary llms. <i>arXiv preprint</i>		
825	<i>arXiv:2305.15717</i> .	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang,	880
		Shayne Longpre, Hwaran Lee, Sangdoon Yun,	881
826	Prakhar Gupta, Chien-Sheng Wu, Wenhao Liu, and	Seongjin Shin, Sungdong Kim, James Thorne, et al.	882
827	Caiming Xiong. 2022. Dialfact: A benchmark for	2023. Prometheus: Inducing fine-grained evalua-	883
828	fact-checking in dialogue. In <i>Proceedings of the</i>	tion capability in language models. <i>arXiv preprint</i>	884
829	<i>60th Annual Meeting of the Association for Compu-</i>	<i>arXiv:2310.08491</i> .	885
830	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages		
831	3785–3801.	Seungone Kim, Juyoung Suk, Shayne Longpre,	886
		Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham	887
832	Or Honovich, Leshem Choshen, Roei Aharoni, Ella	Neubig, Moontae Lee, Kyungjae Lee, and Minjoon	888
833	Neeman, Idan Szpektor, and Omri Abend. 2021.	Seo. 2024. Prometheus 2: An open source language	889
834	Q2: Evaluating factual consistency in knowledge-	model specialized in evaluating other language mod-	890
835	grounded dialogues via question generation and	els. <i>arXiv preprint arXiv:2405.01535</i> .	891
836	question answering. In <i>Proceedings of the 2021</i>		
837	<i>Conference on Empirical Methods in Natural Lan-</i>	Diederik P Kingma and Jimmy Ba. 2015. Adam: A	892
838	<i>guage Processing</i> , pages 7856–7870.	method for stochastic optimization. In <i>International</i>	893
		<i>Conference on Learning Representations</i> .	894
839	Hamish Ivison, Yizhong Wang, Valentina Pyatkin,	Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn	895
840	Nathan Lambert, Matthew Peters, Pradeep Dasigi,	Park, Zae Myung Kim, and Dongyeop Kang.	896
841	Joel Jang, David Wadden, Noah A Smith, Iz Belt-	2023. Benchmarking cognitive biases in large	897
842	agy, et al. 2023. Camels in a changing climate: En-	language models as evaluators. <i>arXiv preprint</i>	898
843	hancing lm adaptation with tulu 2. <i>arXiv preprint</i>	<i>arXiv:2309.17012</i> .	899
844	<i>arXiv:2311.10702</i> .		
845	Shankar Iyer, Nikhil Dandekar, and Kornél Csernai.	Tomasz Korbak, Kejian Shi, Angelica Chen,	900
846	2017. <b>First Quora Dataset release: Question pairs</b> .	Rasika Vinayak Bhalerao, Christopher Buck-	901
		ley, Jason Phang, Samuel R. Bowman, and Ethan	902
847	Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan,	Perez. 2023. <b>Pretraining language models with</b>	903
848	Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun,	<b>human preferences</b> . In <i>Proceedings of the 40th</i>	904
849	Yizhou Wang, and Yaodong Yang. 2023. Beaver-	<i>International Conference on Machine Learning</i> ,	905
850	tails: Towards improved safety alignment of llm via	volume 202 of <i>Proceedings of Machine Learning</i>	906
851	a human-preference dataset. In <i>Thirty-seventh Con-</i>	<i>ference Research</i> , pages 17506–17533. PMLR.	907
852	<i>ference on Neural Information Processing Systems</i>		
853	<i>Datasets and Benchmarks Track</i> .	Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit	908
		Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo.	909
854	Albert Q Jiang, Alexandre Sablayrolles, Antoine	2023a. Longeval: Guidelines for human evaluation	910
855	Roux, Arthur Mensch, Blanche Savary, Chris	of faithfulness in long-form summarization. In <i>Pro-</i>	911
856	Bamford, Devendra Singh Chaplot, Diego de las	<i>ceedings of the 17th Conference of the European</i>	912
857	Casas, Emma Bou Hanna, Florian Bressand, et al.	<i>Chapter of the Association for Computational Lin-</i>	913
858	2024. Mixtral of experts. <i>arXiv preprint</i>	<i>guistics</i> , pages 1650–1669.	914
859	<i>arXiv:2401.04088</i> .		



915	Kalpesh Krishna, Yapei Chang, John Wieting, and Mohit Iyyer. 2022. Rankgen: Improving text generation with large ranking models. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 199–232.	explainability of question answering systems post-deployment. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 926–937.	972 973 974
920	Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> .	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let’s verify step by step. In <i>The Twelfth International Conference on Learning Representations</i> .	975 976 977 978 979
926	Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2023b. Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. <i>Advances in Neural Information Processing Systems</i> , 36.	Bill Yuchen Lin, Khyathi Chandu, Faeze Brahman, Yuntian Deng, Abhilasha Ravichander, Valentina Pyatkin, Ronan Le Bras, and Yejin Choi. 2024. <b>Wildbench: Benchmarking language models with challenging tasks from real users in the wild</b> .	980 981 982 983 984
931	Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. <i>arXiv preprint arXiv:2403.13787</i> .	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252.	985 986 987 988 989
937	Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, et al. 2021. Datasets: A community library for natural language processing. <i>arXiv preprint arXiv:2109.02846</i> .	Yiqi Liu, Nafise Sadat Moosavi, and Chenghua Lin. 2023a. Llms as narcissistic evaluators: When ego inflates evaluation scores. <i>arXiv preprint arXiv:2311.09766</i> .	990 991 992 993
943	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023a. Halueval: A large-scale hallucination evaluation benchmark for large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6449–6464.	Yixin Liu, Alexander R Fabbri, Jiawen Chen, Yilun Zhao, Simeng Han, Shafiq Joty, Pengfei Liu, Dragomir Radev, Chien-Sheng Wu, and Arman Cohan. 2023b. Benchmarking generation and evaluation capabilities of large language models for instruction controllable summarization. <i>arXiv preprint arXiv:2311.09184</i> .	994 995 996 997 998 999 1000
949	Ruosan Li, Teerth Patel, and Xinya Du. 2023b. Prd: Peer rank and discussion improve large language model based evaluations. <i>arXiv preprint arXiv:2307.02762</i> .	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. <b>The flan collection: Designing data and methods for effective instruction tuning</b> . In <i>Proceedings of the 40th International Conference on Machine Learning (ICML)</i> , volume 202 of <i>Proceedings of Machine Learning Research (PMLR)</i> , pages 22631–22648.	1001 1002 1003 1004 1005 1006 1007 1008 1009
953	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023c. Alpacaeval: An automatic evaluator of instruction-following models. <a href="https://github.com/tatsu-lab/alpaca_eval">https://github.com/tatsu-lab/alpaca_eval</a> .	Mounica Maddela, Yao Dou, David Heineman, and Wei Xu. 2023. Lens: A learnable evaluation metric for text simplification. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16383–16408.	1010 1011 1012 1013 1014 1015
958	Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022a. <b>Competition-level code generation with alphacode</b> . <i>Science</i> , 378(6624):1092–1097.	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9004–9017.	1016 1017 1018 1019 1020 1021
969	Zichao Li, Prakhar Sharma, Xing Han Lu, Jackie Chi Kit Cheung, and Siva Reddy. 2022b. Using interactive feedback to improve the accuracy and	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. <i>arXiv preprint arXiv:2005.00661</i> .	1022 1023 1024 1025
		AI Meta. 2024. <b>Introducing meta llama 3: The most capable openly available llm to date</b> . <i>Meta AI</i> .	1026 1027



1028	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis,	Arjun Panickssery, Samuel R. Bowman, and Shi Feng.	1084
1029	Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettle-	2024. Llm evaluators recognize and favor their own	1085
1030	moyer, and Hannaneh Hajishirzi. 2023. Factscore:	generations. <i>arXiv preprint arXiv:2404.13076</i> .	1086
1031	Fine-grained atomic evaluation of factual precision		
1032	in long form text generation. In <i>Proceedings of the</i>	Zarana Parekh, Jason Baldridge, Daniel Cer, Austin	1087
1033	<i>2023 Conference on Empirical Methods in Natural</i>	Waters, and Yinfei Yang. 2021. Crisscrossed cap-	1088
1034	<i>Language Processing</i> , pages 12076–12100.	tions: Extended intramodal and intermodal seman-	1089
		tic similarity judgments for ms-coco. In <i>Proceed-</i>	1090
1035	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and	<i>ings of the 16th Conference of the European Chap-</i>	1091
1036	Hannaneh Hajishirzi. 2022. <a href="#">Cross-task generaliza-</a>	<i>ter of the Association for Computational Linguistics:</i>	1092
1037	<a href="#">tion via natural language crowdsourcing instructions</a> .	<i>Main Volume</i> , pages 2855–2870.	1093
1038	In <i>Proceedings of the 60th Annual Meeting of the</i>		
1039	<i>Association for Computational Linguistics (ACL)</i> ,	Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers,	1094
1040	pages 3470–3487.	John Thickstun, Sean Welleck, Yejin Choi, and Zaid	1095
		Harchaoui. 2021. Mauve: Measuring the gap be-	1096
1041	Seungjun Moon, Yongho Song, Hyungjoo Chae,	tween neural text and human text using divergence	1097
1042	Dongjin Kang, Taeyoon Kwon, Kai Tzu-iunn Ong,	frontiers. <i>Advances in Neural Information Process-</i>	1098
1043	Seung-won Hwang, and Jinyoung Yeo. 2023. Cof-	<i>ing Systems</i> , 34:4816–4828.	1099
1044	fee: Boost your code llms by fixing bugs with feed-		
1045	back. <i>arXiv preprint arXiv:2311.07215</i> .	Yiwei Qin, Weizhe Yuan, Graham Neubig, and Pengfei	1100
		Liu. 2023. T5score: Discriminative fine-tuning of	1101
1046	Niklas Muennighoff, Qian Liu, Armel Randy Ze-	generative evaluation metrics. In <i>Findings of the</i>	1102
1047	baze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo,	<i>Association for Computational Linguistics: EMNLP</i>	1103
1048	Swayam Singh, Xiangru Tang, Leandro Von Werra,	2023, pages 15185–15202.	1104
1049	and Shayne Longpre. 2024. Octopack: Instruction		
1050	tuning code large language models. In <i>The Twelfth</i>	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	1105
1051	<i>International Conference on Learning Representa-</i>	pher D Manning, Stefano Ermon, and Chelsea Finn.	1106
1052	<i>tions</i> .	2023. Direct preference optimization: Your lan-	1107
		guage model is secretly a reward model. <i>Advances</i>	1108
1053	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu,	<i>in Neural Information Processing Systems</i> , 36.	1109
1054	Long Ouyang, Christina Kim, Christopher Hesse,		
1055	Shantanu Jain, Vineet Kosaraju, William Saunders,	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	1110
1056	et al. 2021. Webgpt: Browser-assisted question-	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	1111
1057	answering with human feedback. <i>arXiv preprint</i>	Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the lim-</a>	1112
1058	<i>arXiv:2112.09332</i> .	<a href="#">its of transfer learning with a unified text-to-text</a>	1113
		<a href="#">transformer</a> . <i>Journal of Machine Learning Research</i>	1114
1059	Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu,	( <i>JMLR</i> ), 21(140):1–67.	1115
1060	Huan Wang, Yingbo Zhou, Silvio Savarese, and		
1061	Caiming Xiong. 2023. Codegen: An open large lan-	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	1116
1062	guage model for code with multi-turn program syn-	Lavie. 2020. Comet: A neural framework for mt	1117
1063	thesis. In <i>The Eleventh International Conference on</i>	evaluation. In <i>Proceedings of the 2020 Conference</i>	1118
1064	<i>Learning Representations</i> .	<i>on Empirical Methods in Natural Language Process-</i>	1119
		<i>ing (EMNLP)</i> , pages 2685–2702.	1120
1065	OpenAI. 2024a. <a href="#">GPT-3.5 Turbo</a> .		
1066	OpenAI. 2024b. <a href="#">Hello GPT-4o</a> .	Machel Reid, Nikolay Savinov, Denis Teplyashin,	1121
		Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste	1122
1067	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan	1123
1068	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	Firat, Julian Schrittwieser, et al. 2024. Gemini	1124
1069	Sandhini Agarwal, Katarina Slama, Alex Ray, John	1.5: Unlocking multimodal understanding across	1125
1070	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	millions of tokens of context. <i>arXiv preprint</i>	1126
1071	Maddie Simens, Amanda Askell, Peter Welinder,	<i>arXiv:2403.05530</i> .	1127
1072	Paul F Christiano, Jan Leike, and Ryan Lowe. 2022.		
1073	<a href="#">Training language models to follow instructions</a>	Adam Roberts, Hyung Won Chung, Gaurav Mishra,	1128
1074	<a href="#">with human feedback</a> . In <i>Advances in Neural Infor-</i>	Anselm Levskaya, James Bradbury, Daniel Andor,	1129
1075	<i>mation Processing Systems (NeurIPS)</i> , volume 35,	Sharan Narang, Brian Lester, Colin Gaffney, Afroz	1130
1076	pages 27730–27744.	Mohiuddin, et al. 2023. Scaling up models and data	1131
		with t5x and seqio. <i>Journal of Machine Learning</i>	1132
		<i>Research</i> , 24(377):1–8.	1133
1077	Artidoro Pagnoni, Vidhisha Balachandran, and Yulia		
1078	Tsvetkov. 2021. Understanding factuality in abstrac-	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	1134
1079	tive summarization with frank: A benchmark for fac-	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	1135
1080	tuality metrics. In <i>Proceedings of the 2021 Confer-</i>	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja,	1136
1081	<i>ence of the North American Chapter of the Associ-</i>	et al. 2022. <a href="#">Multitask prompted training enables</a>	1137
1082	<i>ation for Computational Linguistics: Human Lan-</i>	<a href="#">zero-shot task generalization</a> . <i>Proceedings of the</i>	1138
1083	<i>guage Technologies</i> , pages 4812–4829.	<i>10th International Conference on Learning Repre-</i>	1139
		<i>sentations (ICLR)</i> .	1140

1141	Tal Schuster, Adam Fisch, and Regina Barzilay. 2021.	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin	1198
1142	Get your vitamin c! robust fact verification with con-	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-	1199
1143	trastive evidence. In <i>Proceedings of the 2021 Con-</i>	drew M Dai, and Quoc V Le. 2022. <a href="#">Finetuned lan-</a>	1200
1144	<i>ference of the North American Chapter of the Asso-</i>	<a href="#">guage models are zero-shot learners</a> . <i>Proceedings of</i>	1201
1145	<i>ciation for Computational Linguistics: Human Lan-</i>	<i>the 10th International Conference on Learning Rep-</i>	1202
1146	<i>guage Technologies</i> , pages 624–643.	<i>resentations (ICLR)</i> .	1203
1147	Thibault Sellam, Dipanjan Das, and Ankur Parikh.	Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu,	1204
1148	2020. <a href="#">BLEURT: Learning robust metrics for text</a>	Nathan Hu, Dustin Tran, Daiyi Peng, Ruibo Liu,	1205
1149	<a href="#">generation</a> . In <i>Proceedings of the 58th Annual Meet-</i>	Da Huang, Cosmo Du, et al. 2024. Long-form fac-	1206
1150	<i>ing of the Association for Computational Linguistics</i> ,	tuality in large language models. <i>arXiv preprint</i>	1207
1151	pages 7881–7892, Online. Association for Computa-	<i>arXiv:2403.18802</i> .	1208
1152	tional Linguistics.	Adina Williams, Nikita Nangia, and Samuel Bowman.	1209
1153	Ondrej Skopec, Rahul Aralikkatte, Sian Gooding, and	2018. A broad-coverage challenge corpus for sen-	1210
1154	Victor Cărbune. 2023. Towards better evaluation of	tence understanding through inference. In <i>Proceed-</i>	1211
1155	instruction-following: A case-study in summariza-	<i>ings of the 2018 Conference of the North American</i>	1212
1156	tion. In <i>Proceedings of the 27th Conference on Com-</i>	<i>Chapter of the Association for Computational Lin-</i>	1213
1157	<i>putational Natural Language Learning (CoNLL)</i> ,	<i>guistics: Human Language Technologies, Volume 1</i>	1214
1158	pages 221–237.	<i>(Long Papers)</i> , pages 1112–1122.	1215
1159	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel	Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum,	1216
1160	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,	Cheng Niu, Randy Zhong, Juntong Song, and Tong	1217
1161	Dario Amodei, and Paul F Christiano. 2020. Learn-	Zhang. 2023a. Ragtruth: A hallucination corpus	1218
1162	ing to summarize with human feedback. <i>Advances</i>	for developing trustworthy retrieval-augmented lan-	1219
1163	<i>in Neural Information Processing Systems</i> , 33:3008–	guage models. <i>arXiv preprint arXiv:2401.00396</i> .	1220
1164	3021.	Zequ Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane	1221
1165	Yixuan Su and Jialu Xu. 2022. An empirical	Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari	1222
1166	study on contrastive search and contrastive decod-	Ostendorf, and Hannaneh Hajishirzi. 2023b. Fine-	1223
1167	ing for open-ended text generation. <i>arXiv preprint</i>	grained human feedback gives better rewards for lan-	1224
1168	<i>arXiv:2211.10797</i> .	guage model training. In <i>Thirty-seventh Conference</i>	1225
1169	Liyang Tang, Philippe Laban, and Greg Durrett.	<i>on Neural Information Processing Systems</i> .	1226
1170	2024. Minicheck: Efficient fact-checking of	Fangyuan Xu, Junyi Jessy Li, and Eunsol Choi. 2022.	1227
1171	llms on grounding documents. <i>arXiv preprint</i>	How do we answer complex questions: Discourse	1228
1172	<i>arXiv:2404.10774</i> .	structure of long-form answers. In <i>Proceedings</i>	1229
1173	Brian Thompson and Matt Post. 2020. Automatic ma-	<i>of the 60th Annual Meeting of the Association for</i>	1230
1174	chine translation evaluation in many languages via	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	1231
1175	zero-shot paraphrasing. In <i>Proceedings of the 2020</i>	pages 3556–3572.	1232
1176	<i>Conference on Empirical Methods in Natural Lan-</i>	Fangyuan Xu, Yixiao Song, Mohit Iyyer, and Eunsol	1233
1177	<i>guage Processing (EMNLP)</i> , pages 90–121.	Choi. 2023a. A critical evaluation of evaluations	1234
1178	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant,	for long-form question answering. In <i>Proceedings</i>	1235
1179	Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung,	<i>of the 61st Annual Meeting of the Association for</i>	1236
1180	Denny Zhou, Quoc Le, et al. 2023. <a href="#">Freshllms: Re-</a>	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	1237
1181	<a href="#">freshing large language models with search engine</a>	pages 3225–3245.	1238
1182	<a href="#">augmentation</a> . <i>arXiv preprint arXiv:2310.03214</i> .	Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao	1239
1183	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020.	Song, Markus Freitag, William Wang, and Lei Li.	1240
1184	Asking and answering questions to evaluate the fac-	2023b. Instructscore: Towards explainable text gen-	1241
1185	tual consistency of summaries. In <i>Proceedings of</i>	eration evaluation with automatic feedback. In <i>Pro-</i>	1242
1186	<i>ceedings of the 58th Annual Meeting of the Association for</i>	<i>ceedings of the 2023 Conference on Empirical Meth-</i>	1243
1187	<i>Computational Linguistics</i> , pages 5008–5020.	<i>ods in Natural Language Processing</i> , pages 5967–	1244
1188	Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams,	5994.	1245
1189	Makesh Narsimhan Sreedhar, Daniel Egert, Olivier	Xinyan Yu, Sewon Min, Luke Zettlemoyer, and Han-	1246
1190	Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan	naneh Hajishirzi. 2023. Crepe: Open-domain ques-	1247
1191	Swope, et al. 2023. Helpsteer: Multi-attribute	tion answering with false presuppositions. In <i>Pro-</i>	1248
1192	helpfulness dataset for steerlm. <i>arXiv preprint</i>	<i>ceedings of the 61st Annual Meeting of the Associa-</i>	1249
1193	<i>arXiv:2311.09528</i> .	<i>tion for Computational Linguistics (Volume 1: Long</i>	1250
1194	Alex Warstadt, Amanpreet Singh, and Samuel R Bow-	<i>Papers)</i> , pages 10457–10480.	1251
1195	man. 2019. Neural network acceptability judgments.	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.	1252
1196	<i>Transactions of the Association for Computational</i>	<a href="#">Bartscore: Evaluating generated text as text gener-</a>	1253
1197	<i>Linguistics</i> , 7:625–641.	<a href="#">ation</a> . In <i>Advances in Neural Information Process-</i>	1254

1255 *ing Systems*, volume 34, pages 27263–27277. Curran Associates, Inc.

1256

1257 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

1258

1259

1260

1261 Yuan Zhang, Jason Baldridge, and Luheng He. 2019. Paws: Paraphrase adversaries from word scrambling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1298–1308.

1262

1263

1264

1265

1266

1267

1268 Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. 2023. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.

1269

1270

1271

1272 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023a. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.

1273

1274

1275

1276

1277

1278

1279

1280 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 llm-as-a-judge with mt-bench and chatbot arena](#). *Preprint*, arXiv:2306.05685.

1281

1282

1283

1284

1285

1286 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

1287

1288

1289

1290

## Appendix 1291

### A Appendix Tables 1292

#### A.1 List of Training Datasets in FLAME 1293

Please see [Table 6](#). 1294

### B Further analysis on the REWARDBENCH benchmark 1295

In this section, we provide some analysis of issues we found in the REWARDBENCH benchmark, including issues of length bias ([Appendix B.1](#)) and difficulty in hill-climbing ([Appendix B.2](#)). Given these issues, we encourage readers and future efforts to not solely rely on REWARDBENCH for autorater task performance, but instead use a wide variety of evaluation tasks to compare models (such as our evaluation suite in §4). 1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305

#### B.1 Analysis of length and token bias in REWARDBENCH 1306

In [Table 5](#), we present an analysis of length bias in REWARDBENCH. Overall, we find that REWARDBENCH is far from length-balanced in its composition. The Chat-Hard, Coding and Math categories strongly prefer shorter length outputs, while the Chat category has a strong preference towards longer outputs. An adversarial submission which can identify prompt categories may simply choose to prefer the longer or shorter output in REWARDBENCH, and achieve high scores without being an actually strong preference model. 1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318

RewardBench Split	% longer preferred
Chat	79.1%
Chat-Hard	29.6%
Math	6.5%
Coding	35.7%
Safety	41.9%

Table 5: A summary of length bias in REWARDBENCH. Overall, we find that four out of five categories in REWARDBENCH have a strong preference towards longer or shorter outputs.

Besides length bias, we also found some issues of token bias in the Math and Safety splits of REWARDBENCH. In Safety, the preferred side had a strong preference for phrases like “I’m sorry”, which are indicative of hedged responses. In 28% pairs, only the preferred response contained the word “sorry”. We found similar issues in the Math split, with tokens “i”, “can”, “need”, “to”, “find” largely only appearing in the rejected response. 1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327

1328 Given these findings, our recommendation for  
1329 future autorater / reward model development is  
1330 **not just rely on REWARDBENCH performance, but**  
1331 **instead to evaluate a wide variety of autorater**  
1332 **tasks** (such as our evaluation suite in §4).

## 1333 **B.2 Discussion on hill-climbing on** 1334 **REWARDBENCH**

1335 In early experiments, we found it very difficult to  
1336 hill-climb on REWARDBENCH due to the absence of a  
1337 development set in Lambert et al. (2024). It was not  
1338 possible to construct a “proxy” development set for  
1339 many categories of REWARDBENCH, since Lambert  
1340 et al. (2024) had fully utilized the original set while  
1341 constructing them (such as LLMBAR in CHATHARD).  
1342 In early experiments, we saw poor correlation be-  
1343 tween REWARDBENCH performance and performance  
1344 on other held-out tasks between our model variants.  
1345 This prompted us to hill-climb directly on REWARD-  
1346 BENCH as a proxy, as described in §3.4. To confirm  
1347 we have not overfit on REWARDBENCH, we evalu-  
1348 ate both FLAME and FLAME-RM on several other  
1349 held-out tasks besides REWARDBENCH, to confirm  
1350 the general purpose utility of our models across a  
1351 variety of tasks.



Capability	Dataset	Source	Output Format
Attribution / Factuality	ESNLI	Camburu et al. (2018)	Classification, Generative
	MNLI	Williams et al. (2018)	Classification
	VitaminC	Schuster et al. (2021)	Classification
	Sentence Similarity - CxC	Parekh et al. (2021)	Pointwise
	Sentence Similarity - STSB	Cer et al. (2017)	Pointwise
	MultiPIT	Dou et al. (2022b)	Classification
	QQP	Iyer et al. (2017)	Classification
	PAWS Paraphrasing	Zhang et al. (2019)	Classification
	FaithDial	Dziri et al. (2022a)	Classification
	MOCHA	Chen et al. (2020)	Pointwise
	DialFact	Gupta et al. (2022)	Classification
	RAGTruth	Wu et al. (2023a)	Classification
	FActScore	Min et al. (2023)	Classification
	FRANK	Pagnoni et al. (2021)	Classification
	BEGIN	Dziri et al. (2022b)	Classification
	XSUM-Faithful	Maynez et al. (2020)	Generative
	HaluEval	Li et al. (2023a)	Classification
	QAGS	Wang et al. (2020)	Classification
	WikiBio Hallucinations	Manakul et al. (2023)	Pointwise
	Q2	Honovich et al. (2021)	Classification
General Text Quality	GENIE	Khashabi et al. (2021)	Pointwise, Pairwise, Gen.
	InstruSum	Liu et al. (2023b)	Pairwise, Classification
	RiSum	Skopek et al. (2023)	Pointwise, Classification
	Stanford SHP	Ethayarajh et al. (2022)	Pairwise
	BeaverTails Helpful	Ji et al. (2023)	Pairwise
	HH RLHF Helpful	Bai et al. (2022)	Pairwise
	Summary Feedback Comparisons	Stiennon et al. (2020)	Pairwise, Pointwise
	SEAHORSE	Clark et al. (2023)	Classification
	Scarecrow	Dou et al. (2022a)	Classification, Generative
	SummaEval	Fabbri et al. (2021)	Pointwise
	LMSys Chatbot Arena (english)	Zheng et al. (2023a)	Pairwise
	FeedbackQA	Li et al. (2022b)	Pointwise, Generative
	WebGPT	Nakano et al. (2021)	Pairwise, Generative
	Fine-grained RLHF	Wu et al. (2023b)	Pairwise, Classification
	LENS	Maddela et al. (2023)	Pointwise
	MAUVE - Human Eval	Pillutla et al. (2021)	Pairwise
	CoLA	Warstadt et al. (2019)	Classification
	CREPE	Yu et al. (2023)	Classification, Generative
	PRD-Vicuna	Li et al. (2023b)	Pairwise
	Hurdles LFQA	Krishna et al. (2021)	Pairwise
	Validity LFQA	Xu et al. (2022)	Classification, Generative
	News Summarization Evaluation	Goyal et al. (2022)	Pairwise
	Helpful Steer (training split)	(Wang et al., 2023)	Pointwise
Safety	BeaverTails Classify	Ji et al. (2023)	Classification
	BeaverTails Harmless	Ji et al. (2023)	Pairwise
	HH RLHF Harmless	Bai et al. (2022)	Pairwise
	HH RLHF Red Teaming	Bai et al. (2022)	Pointwise
Coding	Code Contests	Li et al. (2022a)	Pairwise
	COFFEE	Moon et al. (2023)	Pairwise
	CommitPack	Muennighoff et al. (2024)	Pairwise
	CommitPack - Bugs	Muennighoff et al. (2024)	Pairwise
Math Reasoning	PRM800K preference	Lightman et al. (2024)	Pairwise
Instruction Tuning	TULU V2	Iverson et al. (2023)	Generative
	PRM800K-IF	Lightman et al. (2024)	Generative
	LIMA	Zhou et al. (2023)	Generative

Table 6: A complete list of datasets which were used to train FLAME, along with their output format and capability categorization. From many source datasets, we derived multiple training dataset tasks (for example, by splitting the pointwise and pairwise ratings for the the same set of responses).