EXPLAINING LENGTH BIAS IN LLM-BASED PREFER ENCE EVALUATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

The use of large language models (LLMs) as judges, particularly in preference comparisons has become widespread, but this reveals a notable bias towards longer responses, undermining the reliability of such evaluations. To better understand such bias, we propose to decompose the preference evaluation metric, specifically the *win rate*, into two key components: *desirability* and *information mass*, where the former is length-independent and related to trustworthiness such as correctness, toxicity, and consistency, and the latter is length-dependent and represents the amount of information in the response. We empirically demonstrated the decomposition through controlled experiments and found that response length impacts evaluations by influencing information mass. To derive a reliable evaluation metric that assesses content quality without being confounded by response length, we propose AdapAlpaca, a simple yet effective adjustment to win rate measurement. Specifically, AdapAlpaca ensures a fair comparison of response quality by aligning the lengths of reference and test model responses under equivalent length intervals.

025 026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

027 028 029

As LLMs are increasingly deployed across various domains of artificial intelligence, from natural language processing to complex decision-making systems (Wu et al., 2023; Li et al., 2023a; Rao 031 et al., 2023; Song et al., 2023), ensuring their performance, reliability, and fairness has become a 032 critical challenge (Louis & Nenkova, 2013; Wang et al., 2023b). LLM-based auto-evaluators have 033 emerged as a crucial tool in this context, offering a cost-effective and scalable alternative to labor-034 intensive human evaluations (Chen et al., 2023; Li et al., 2024a;b; Dubois et al., 2024b). Despite their advantages, these automated systems are not without their shortcomings, particularly concern-035 ing the introduction and perpetuation of biases (Li et al., 2023c; Zheng et al., 2023; Koo et al., 036 2023; Wang et al., 2023a; Wu & Aji, 2023). One of the important biases observed in LLM-based 037 evaluations is the preference for longer textual responses. Previous empirical studies have explored a strong correlation between the length of response and its perceived quality represented by win rate (Zhao et al., 2024; Dubois et al., 2024a; Ivison et al., 2023). However, it is not reasonable to 040 simply attribute the preference to length since length is only the surface factor for the quality of a 041 sentence. Therefore, in this work, we investigate the following question: what are the major factors 042 contributing to the win rate? 043

To solve this problem, we propose a new framework that decomposes the *quality* of a response, as 044 measured by its win rate in pairwise comparisons, into two distinct components: (1) desirability, 045 which is independent of length and reflects the trustworthiness of the response, encompassing fac-046 tors such as correctness, toxicity, and consistency; and (2) information mass, which is dependent 047 on length and represents the amount of information in the response, measurable through conditional 048 entropy. We validate our hypothesis by testing win rates in two different scenarios: (i) comparing normal responses with those differing in desirability (e.g., Logical to be desired and Biased not desired), and (ii) comparing normal responses with concise and detailed responses, which vary in in-051 formation mass. Our experiments demonstrate that responses with negative desirability significantly decrease the win rate, whereas information mass, when not negatively influenced by desirability, 052 is positively correlated with the win rate, thus confirming the effectiveness of our metric. Following this finding, we design a new prompt called "Quality Enhancement" to improve information



Figure 1: Comparison between AlpacaEval and AdapAlpaca (Ours). In AlpacaEval, the reference answer has a fixed length, regardless of the length of the test model's answer. In contrast, AdapAlpaca dynamically selects a reference answer that matches the length of the test model's answer.

mass with positive desirability. This prompt enables GPT-4 to achieve state-of-the-art results on
 AlpacaEval, increasing the win rate from 50.00% to 70.16%.

072 Through our decomposition of the quality of a response, we observe that response length impacts 073 evaluations primarily by influencing information mass. However, a reliable evaluation metric should 074 assess content quality without being confounded by extraneous factors such as response length (Koo 075 et al., 2023; Ye et al., 2024; Dubois et al., 2024a), we further propose AdapAlpaca, a benchmark 076 designed to improve evaluation fairness. By ensuring that responses are compared at the same 077 length intervals, AdapAlpaca effectively mitigates length bias, enabling accurate content quality 078 assessments (see Figure 1). With AdapAlpaca, we further analyze length bias in Direct Preference 079 Optimization (DPO) Rafailov et al. (2023) to examine the findings in prior work (Gu et al., 2024; Ivison et al., 2023; Liu et al., 2024) that DPO lengthens model responses. Specifically, we test TÜLU2 (Ivison et al., 2023) and TÜLU2-dpo models at 7B, 13B, and 70B scales on AlpacaEval 081 and AdapAlpaca. Our results indicate that DPO leads to higher human preference, but this gain is amplified by response length, with AlpacaEval showing higher win rates gain than AdapAlpaca. 083 Our major findings and contributions are as follows: 084

- We propose a novel interpretation of win rate, emphasizing desirability and information mass, offering a more precise LLM performance measure. Based on this interpretation, we develop the "Quality Enhancement" prompt, which improves win rates by boosting information mass with positive desirability. This prompt improves win rates across multiple LLMs, with average increases of 23.44% for GPT-3.5, 16.48% for GPT-4, 22.28% for LLAMA3-70b, and 20.40% for Qwen1.5 72B.
 - To mitigate length bias, we introduce AdapAlpaca, a method that aligns the response lengths of the reference and test model, enabling a fair comparison of desirability and information mass under the same length intervals.
 - Using both AlpacaEval and AdapAlpaca, we analyze the impact of length bias in DPO. Our experiments with TÜLU2 and TÜLU2-dpo models at 7B, 13B, and 70B scales show that DPO leads to higher human preference, but this gain is amplified by response length, with AlpacaEval showing higher win rates gain than AdapAlpaca.
- 098 099 100

101

102

091

092

094

095

096

066

067

068 069

2 RELATED WORK

2.1 REFERENCE-FREE EVALUATION METRICS

Reference-free evaluation metrics have a long history (Louis & Nenkova, 2013), which evaluates the
generated text based on intrinsic properties and coherence with the context. Although they achieve
high accuracy on matching inner-annotator, the achievement suffers from spurious correlations such
as perplexity and length (Durmus et al., 2022). Recently, people have started using a strong model
(e.g., GPT-4) as an evaluator to perform a zero-shot reference-free evaluation on the weak models (Shen et al., 2023; Dubois et al., 2024b; Chen et al., 2023; Hu et al., 2024). However, leveraging

108 a strong model's intrinsic knowledge to perform reference-free evaluation ignores the prompt pref-109 erence of the strong model, for example, the prompt's length.

- 110
- 111 112

2.2 CORRELATION BETWEEN LENGTH AND WIN RATE

Previous research reveals that sentence length will influence the evaluation of trustworthiness. 113 114 Specifically, when using a GPT-4 to represent human preference, it will prefer to choose a long sentence rather than a short sentence (Dubois et al., 2024a; Ivison et al., 2023; Gu et al., 2024; Shen 115 et al., 2023; Koo et al., 2023; Wang et al., 2023a; Wu & Aji, 2023; Dubois et al., 2024b; Chen et al., 116 2023; Hu et al., 2024). Such preference will introduce a length-correlated bias and help the model 117 with long-generation sentences gain a high score on human preference evaluation. Although these 118 approaches show a high correlation to human preference, debiasing such as automated evaluation is 119 highly valuable. (Dubois et al., 2024a) proposes a length-controlled (LC) win rate by removing the 120 length-correlated term in the win rate regression model. The new LC win rate shows an even per-121 formance between concise and verbose input and a higher correlation when compared with human 122 preference.

- 123
- 124
- 125

3 UNDERSTANDING THE MAJOR FACTORS OF WIN RATE

126 To interpret the correlation between length and win rate, we propose a new framework based on 127 quality, which includes desirability (length-independent, related to trustworthiness) and information 128 mass (length-dependent, represented by conditional entropy). We validate our hypothesis through 129 two scenarios: (1) testing the impact of different desirability on win rate with the same informa-130 tion mass, and (2) testing the influence of different information mass on win rate with the same 131 desirability. 132

133 134

3.1 PRELIMINARY

135 **Evaluation protocol.** We utilize the AlpacaEval dataset (Li et al., 2023c) to assess human prefer-136 ences. AlpacaEval is a reference-free evaluation dataset for LLMs, encompassing 805 instructions 137 that reflect human interactions on the Alpaca web demo. To ensure a comprehensive evaluation of human preferences, we extend our testing to additional datasets, including LIMA (Zhou et al., 138 2023), Vicuna (Chiang et al., 2023), Koala (Vu et al., 2023), Wizardlm (Xu et al., 2023), and Self-139 Instruct (Wang et al., 2022), in line with previous studies (Chen et al., 2023; Zhang et al., 2024; Du 140 et al., 2023; Zhao et al., 2024; Li et al., 2023b). 141

142

158

161

Base Models. In our experiments, we follow the setup in the AlpacaEval Leaderboard¹, using 143 the GPT-4 Preview (11/06) as *Baseline* and the *Annotator*. The references to GPT-3.5, LLAMA3-144 70b, and Qwen1.5 72b in the main text denote gpt-3.5-turbo-0125, meta-llama/Meta-Llama-3-70B-145 Instruct², and Qwen/Qwen1.5-72B-Chat³, respectively. Following previous work (Wei et al., 2024), 146 we calculate conditional entropy using the method described in (Von Neumann, 2013). 147

148 Win rate. Assume we have a set of instructions x. We prompt a test model m to generate a response z_m for each instruction. Similarly, we prompt a reference model b (referred to as the 149 150 "baseline" in AlpacaEval) to generate a response z_b for each instruction.⁴ An annotator then evaluates these responses based on their quality and assigns a preference $y \in \{m, b\}$, indicating which 151 model's response is superior. To properly understand the concept of win rate, we first need to define 152 what we mean by response quality: 153

154 **Definition 1** (Response Quality), denoted as $Q_e(z|x)$, quantifies the effectiveness of the model's 155 response z in addressing the given instruction x, as evaluated by an annotator e. Annotator prefer 156 responses with higher quality. 157

¹https://tatsu-lab.github.io/alpaca_eval

²https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct 159

³https://huggingface.co/Qwen/Qwen1.5-72B-Chat 160

 $^{^{4}}$ In this context, m stands for "model" and b denotes "baseline", which in this paper follows the AlpacaEval Leaderboard's use of GPT-4 Preview (11/06).

By leveraging the definition of quality, we can now formulate the win rate as the comparison of sentence quality as follows:

$$WinRate(m,b) = \mathbb{E}_x \left[\mathbbm{1}_{Q_e(z_m|x) > Q_e(z_b|x)} \right], \tag{1}$$

where $\mathbb{1}$ is an indicator function and $\mathbb{1}_{Q_e(z_m|x)>Q_e(z_b|x)}$ represents the preference distribution for each individual. Previous works (Chen et al., 2023; Li et al., 2023c; Dubois et al., 2024a;b) utilize LLMs as zero-shot evaluators due to their exceptional performance on real-world tasks. Our experimental setup adheres to the AlpacaEval Leaderboard ⁵ guidelines, employing the GPT-4 Preview (11/06)⁶ as both the *Baseline b* and the *Annotator e*.



Figure 2: Validation of desirability's impact on quality for GPT-4. The results demonstrate that desirability influences the win rate.

3.2 QUALITY DECOMPOSITION

Before discussing the composition of quality, we first define two key concepts: **desirability** and **information mass**. Desirability reflects the inherent quality attributes of a response that make it reliable and valuable, irrespective of its length, while information mass captures the quantity of information in the response, with longer responses generally containing more content. The definitions of desirability and information mass are as follows:

Definition 2 (*Desirability*), denoted as $D_e(z|x)$, measures the probability that annotator e will accept the response z given an instruction x. It can be influenced by factors such as consistency and toxicity and is independent of response length.

Definition 3 (*Information mass*), denoted as $H_e(z|x)$, measures the amount of information in a response z given an instruction x, as evaluated by annotator e. It is represented by conditional entropy and is directly with response length.

With these definitions in place, we now present our main hypothesis on answer quality, starting withan assumption:

214 215

194

195 196 197

198 199

200

201

202

203 204

205

206

207

⁵https://tatsu-lab.github.io/alpaca_eval

⁶In this paper, unless specified otherwise, GPT-4 refers to GPT-4 Preview (11/06).



Figure 3: Validation of information mass's impact on quality for GPT-4. The results demonstrate that information mass influences the win rate.

Assumption 1 (*Quality Decomposition*). For a given answer z and instruction x, the quality $Q_e(z|x)$ recognized by annotator e can be decomposed as:

$$Q_e(z|x) \propto D_e(z|x) + H_e(z|x), \tag{2}$$

where $D_e(z|x)$ denotes the desirability of the response, and $H_e(z|x)$ represents the information mass.

241 To systematically verify our hypothesis, we conduct two experiments targeting the manipulation of these key components in GPT-4's responses in Section 3.3 and Section 3.4. Additional results with 242 more test and annotator models are provided in Appendix J and Appendix K. 243

3.3 DESIRABILITY INFLUENCES QUALITY

To evaluate the impact of desirability on quality, we design experiments using eight strategies to 247 manipulate response desirability. These strategies include: (1) Origin: No prompt restrictions. (2) 248 **Copy-paste**: Copy GPT-4's response three times. (3) **Biased**: Provide biased responses, favoring 249 certain ideas without justification. (4) Inconsistent: Provide contradictory information to create 250 confusion. Illogical: Give responses based on flawed logic or irrelevant information. (5) Verbose: 251 Provide lengthy responses filled with broad, unrelated details. (6) Toxic: Use offensive language 252 with an aggressive tone. (7) **Relevant**: Provide responses that align with the query. (8) **Logical**: 253 Base responses on sound reasoning and valid arguments. The results are shown in Figure 2. To 254 eliminate the impact of information mass on win rate, we use conditional entropy to represent information mass and ensure the information mass of **Origin** and the other prompts remains as consistent 255 as possible. The entropy values shown in the Figure 2 represent the average conditional entropy of 256 the responses for each prompt. Details for these prompts and relevant implementation are shown in 257 Appendix I and Appendix A.1. First, we observe that although the Copy-paste and Origin prompts 258 maintain identical information mass (as simply replicating text does not increase information), the 259 win rates of **Copy-paste** fall below **Origin** (50%) due to significant consistency impairments. Sec-260 ond, responses generated from negative prompts (i.e., Biased, Inconsistent, Illogical, Verbose, and 261 Toxic) exhibit low desirability, resulting in win rates substantially lower than Origin (50%), despite 262 having similar information mass. Conversely, prompts enhancing desirability (i.e., Consistent and 263 Logical) yield increased win rates compared to Origin. In summary, desirability plays a significant 264 role in determining quality.

265

231

232 233 234

235

236 237 238

239

240

244 245

246

266 3.4 INFORMATION MASS INFLUENCES QUALITY

267

To evaluate the impact of information mass on quality, we designed experiments using three distinct 268 strategies to manipulate the information mass of responses. These strategies include: (1) Origin: 269 No prompt restrictions. (2) Concise: Request brief responses focusing on the most crucial points.

(3) **Detailed**: Request comprehensive responses covering all relevant aspects thoroughly. The cor-responding results are illustrated in Figures 3. Importantly, to isolate the effect of information mass, we ensured that the prompts did not impose any constraints on desirability, ensuring comparability. Details of the prompts and implementation are in Appendix I and Appendix A.1. Our findings indi-cate that information mass significantly affects the win rate without a negative desirability prompt. Specifically, responses with higher information mass, measured by conditional entropy, consistently achieved higher win rates. Thus, we observe the following relationship: Detailed > Origin > Concise. These results confirm that information mass is a crucial factor influencing the quality of responses.

Table 1: The content of the "Quality Enhancement" prompt, designed to elevate both the information mass and desirability of responses, thereby enhancing win rates. Keywords such as "relevant" and "logical" are used to enhance desirability, while "detailed" is used to boost information mass.

Quality EnhancementYou are an expert assistant, delve deeply into the core of the topic, providing a richly
detailed response that explores all its dimensions. Ensure each part of your response is
relevant to the query in a logical manner. Your response should provide comprehensive
information and thoroughly cover all relevant aspects with accuracy and depth.



Figure 4: Correlation between information mass and word count for responses of GPT-4. As the word count increases, the information mass also increases.

3.5 QUALITY ENHANCEMENT PROMPT

Our decomposition reveals that responses with good desirability and higher information mass are generally more favored. Building on this insight, we propose the 'Quality Enhancement' prompt (Ta-ble 1), designed to improve both desirability and information mass, thereby increasing win rates. The keywords "relevant" and "logical" are used to enhance desirability, while "detailed" is used to boost information mass. Their effectiveness is validated in Section 3.2. We evaluated this prompt across multiple models, including GPT-3.5, GPT-4, LLAMA3-70b, and Qwen1.5 72B. The results, summarized in Table 2, with benchmarks such as LIMA, Vicuna, Koala, Wizardlm, and Self-Instruct. The consistent improvement in win rates across all tested models underscores the critical role of response quality in LLM evaluation.

326

327 328

340 341

342 343

Table 2: Win rates with and without the "Quality Enhancement" prompt, along with the corresponding win rate gains (WR Gain). "WR Gain" represents the increase in win rate due to the use of the "Ouality Enhancement".

Models	Methods	AlpacaEval	LIMA	Koala	Self-instruct	Vicuna	Wizardlm	Avg.
	w/o Quality Enhancement	15.47	9.67	11.39	21.46	8.75	16.82	13.93
GPT-3.5	with Quality Enhancement	29.89	36.53	40.34	45.93	35.88	35.62	37.36
	WR Gain	14.42	26.86	28.95	24.47	27.13	18.80	23.44
	w/o Quality Enhancement	50.00	50.00	50.00	50.00	50.00	50.00	50.00
GPT-4	with Quality Enhancement	70.16	65.84	58.90	67.06	73.13	63.76	66.48
	WR Gain	20.16	15.84	8.90	17.06	23.13	13.76	16.48
	w/o Quality Enhancement	34.32	36.63	40.12	39.70	36.74	36.99	37.81
LLAMA3-70b	with Quality Enhancement	56.50	60.39	61.30	64.81	63.49	51.70	59.70
	WR Gain	22.18	23.76	21.18	25.11	26.75	14.71	22.28
	w/o Quality Enhancement	28.27	28.40	35.25	33.81	33.70	31.80	32.67
Qwen1.5 72b	with Quality Enhancement	48.87	53.34	55.40	52.43	56.49	47.13	52.28
	WR Gain	20.60	24.94	20.15	18.62	22.79	15.33	20.40

4 ADAPTIVE ALPACAEVAL

4.1 MOTIVATION

344 345

346 Here, we analyze the phenomenon observed in prior works (Dubois et al., 2024a; Chen et al., 2023; 347 Dubois et al., 2024b), which highlights a positive correlation between response length and win rate. 348 Intuitively, longer responses tend to encompass more information. To rigorously quantify this rela-349 tionship, we use conditional entropy as information mass in a response z given an instruction x. This 350 analysis is conducted without constraints on response desirability, ensuring the correlation between 351 length and information mass remains independent of desirability factors. As shown in Figure 4, our 352 analysis demonstrates a clear trend: as the length of a response increases, the information mass also 353 grows. By integrating this observation with the findings from Section 3.2, we conclude that the primary mechanism through which length affects win rate is its contribution to the overall information 354 mass. 355

356 Adaptive AlpacaEval (AdapAlpac) is based on the premise that a reliable evaluation metric should 357 not only assess the content quality but also ensure that the assessment is not confounded by extrane-358 ous factors such as the length of the response. Central to this approach is the concept of information mass, which is inherently dependent on response length and can be quantified using conditional en-359 tropy. Our primary aim is to mitigate scenarios where merely extending the length of a response ar-360 tificially inflates its conditional entropy and, thus, its perceived quality by annotators. This approach 361 involves dynamically adjusting the evaluation criteria based on response length, thereby providing a 362 more equitable and accurate measure of a model's performance. 363

364 365

4.2 DATASET GENERATION

366 367

To support the development of Adaptive AlpacaEval, we first generate a diverse dataset using a 368 modified prompting strategy with GPT-4, designed to produce responses within specific word count 369 ranges. Specifically, we analyzed the word count distribution within the AlpacaEval dataset, ob-370 serving that responses predominantly fall within the 0-1000 word range. This range was chosen 371 to encompass the full spectrum of response lengths present in the original AlpacaEval dataset, 372 ensuring comprehensive evaluation coverage. To systematically explore this range, we divided 373 it into five equal segments, each representing a distinct dataset: AdapAlpaca-200: 0-200 words, 374 AdapAlpaca-400: 200-400 words, AdapAlpaca-600: 400-600 words, AdapAlpaca-800: 600-800 375 words, AdapAlpaca-1000: 800-1000 words. Each segment is populated by generating responses using the dataset generation prompt, with GPT-4 configured to produce responses that strictly conform 376 to the specified word counts. The data generation prompt and additional details for AdapAlpaca can 377 be found in Appendix H.

Table 3: Comparison of five quantitative metrics related to quality: Vocabulary Size, Win Rate relative to AlpacaEval (AlpacaWR), Entropy, Inter-sample N-gram Frequency (INGF), and Word Counts.

Interval	Vocabulary Size		AlpacaWR		Entropy		INGE	Word Counts	
	All	Ans Avg.	WR	LCWR	All	Ans Avg.	ntor	Word Counts	
AlpacAns Origin	38474	47.79	50.00	50.00	408.83	0.5686	7376.92	363.85	
AdapAlpaca-200	22612	28.08	20.73	43.81	363.55	0.5056	1618.69	145.72	
AdapAlpaca-400	36943	45.89	47.34	47.40	414.39	0.5763	6003.87	355.20	
AdapAlpaca-600	47691	59.24	62.58	50.97	434.77	0.6046	9086.01	540.95	
AdapAlpaca-800	55362	68.77	71.20	54.31	447.48	0.6223	10320.11	708.36	
AdapAlpaca-1000	66095	82.10	66.98	36.24	456.32	0.6346	10981.84	913.44	

4.3 ANALYSIS OF THE GENERATED DATA

The analysis is structured to quantify each dataset's basic characteristics, followed by a comparative assessment to identify any significant differences attributable to the varying response lengths. Table 3 presents a comprehensive overview, providing a snapshot of the informational content across different datasets. Specifically, it includes vocabulary size, inter-sample N-gram Frequency (INGF) (Mishra et al., 2020), word counts of the generated dataset, win rate, length-controlled win rate, and entropy for AlpacaEval-Origin and AdapAlpaca-200, AdapAlpaca-400, AdapAlpaca-600, AdapAlpaca-800, and AdapAlpaca-1000. Our findings indicate the following: 1) Longer responses generally exhibit higher vocabulary sizes and word counts, suggesting a richer linguistic structure. 2) The INGF metric reveals that while longer responses tend to include more common N-grams, there is significant variability in the types of N-grams used, indicating a creative and diverse use of language. 3) Under Win Rate (WR) metrics, longer responses disproportionately receive higher preference scores due to their higher information mass. However, applying the length-controlled win rate (LCWR) significantly mitigates this bias, leading to a more balanced distribution of scores across different response lengths. This analysis aims to ascertain whether this phenomenon is intrin-sic to the response quality or merely a byproduct of increased length. Our results demonstrate that although longer responses generally possess higher information mass, the quality of information, as measured by win rate, does not necessarily increase proportionally. Excessively lengthy responses can result in a decline in desirability, such as reduced consistency. For instance, in Table 3, the win rate of AdapAlpaca-1000 is lower than that of AdapAlpaca-800.

Table 4: The subscripts in the LCWR and WR columns indicate the differences between these metrics and the corresponding Human WR. A larger absolute value denotes a greater disparity between the annotator's evaluation and Human Preference. "LLM Response" denotes different responses to AlpacaEval questions for GPT-4, with detailed content available in Section 3.2.

LLM Response		AlpacaEval	AdapAlpaca		
	Human	LCWR	WR	Human	WR
Concise	10.81	$35.16_{+24.35}$	$15.96_{+5.15}$	29.56	$28.44_{-1.12}$
Detailed	61.61	$54.13_{-7.48}$	$65.83_{\pm 4.22}$	56.02	$55.36_{-0.78}$
Quality Enhancement	66.70	$49.37_{-17.33}$	$70.16_{\pm 3.46}$	58.88	$57.81_{\pm 1.07}$

4.4 **RESULT OF HUMAN EVALUATION**

Table 4 presents the results of the human study, with details provided in Appendix B and you can find more results in Appendix G. First, we test the results of concise, detail, and quality enhancement (descriptions provided in Section 3 and 3.5) using AlpacaEval, followed by AdapAlpaca. From the gap values between LCWR and human evaluations, we observe significant misalignments, indicating inherent problems with the LCWR metric. In contrast, the win rate calculated using AdapAlpaca closely aligns with the human results, showing an average difference of 0.99% (1.12% + 1.07% + 0.78% / 3). Additionally, we find that the difference between human evaluation and WR decreases as

Table 5: Win rate and response length comparison for TÜLU2 models (7B, 13B, and 70B) on
AlpacaEval and AdapAlpaca. The results indicate that while DPO increases response length and
improves win rate, the win rate gain is further amplified by the response length, leading to higher
performance in AlpacaEval compared to AdapAlpaca.

Size	Model	Winra	Avg Length	
5120	Wieder	AlpacaEval	AdapAlpaca	ing. Bongui
	TÜLU 2	3.60	5.84	203.60
7B	TÜLU 2+DPO	8.33	9.04	282.92
	Gain from DPO	4.73	3.20	-
	TÜLU 2	4.35	8.07	192.58
13B	TÜLU 2+DPO	10.82	13.17	276.96
	Gain from DPO	6.47	5.10	-
70B	TÜLU 2	7.34	10.94	184.26
	TÜLU 2+DPO	15.67	17.90	267.23
	Gain from DPO	8.33	6.96	-

452 the quality of responses improves (from concise to detailed to Quality Enhancement). This suggests 453 that as response quality increases, the preferences of annotators and human evaluators converge. 454 Moreover, we found that the smallest difference in win rate between GPT-4 and human evaluations 455 occurs when using the "Quality Enhancement" prompt, which has the highest levels of desirability 456 and information mass. This further underscores the importance of enhancing both desirability and 457 information mass in model responses. Overall, while both AdapAlpaca and LCWR aim to mitigate length bias in evaluating human preferences, their approaches differ fundamentally. AdapAlpaca 458 eliminates length bias from the outset, whereas LCWR attempts to correct for length bias after it 459 has already influenced the evaluation. The inherent issue with LCWR is that length significantly 460 impacts human preference, and adjusting for length retrospectively is not a reliable approach. 461

462 4.5 DPO AND ITS LENGTH BIAS

463 Previous work (Gu et al., 2024; Ivison et al., 2023) has shown that DPO Rafailov et al. (2023) tends 464 to make model responses longer, raising a natural question: Does the increase in human preference 465 brought by DPO partly stem from the length of the responses? In other words, does DPO generate 466 longer replies, thereby increasing their win rate? To investigate this issue, we conducted tests using 467 the widely-used TULU2 (Ivison et al., 2023) series models. As shown in Table 5, we tested the 468 models at 7B, 13B, and 70B scales on both AlpacaEval and AdapAlpaca to measure their win rates and corresponding response lengths. The results from AlpacaEval and AdapAlpaca indicate 469 that while DPO does lead to longer model responses, it enhances the model's human preference 470 capability (as evidenced by the increased win rate in AdapAlpaca). However, this gain is amplified 471 by the response length (as the win rate in AlpacaEval is higher than in AdapAlpaca). Additionally, 472 we found that all models have higher win rates on AdapAlpaca compared to AlpacaEval. This is 473 because the responses from GPT-4 (1106) on AlpacaEval are longer (363 words, see Appendix 4.3), 474 which unfairly amplifies the capabilities of GPT-4 due to its length. These results emphasize the 475 need for length control in evaluations to reflect true model performance.

476 477 478

5 CONCLUSION

In this paper, we identify and address the critical issue of length bias in LLM-based preference
evaluations, which undermines the reliability of win rate metrics. By decomposing win rate into
desirability and information mass, we offer a nuanced understanding of response quality. Our proposed framework, AdapAlpaca, effectively mitigates length bias by dynamically adjusting reference
answer lengths to match test model responses, ensuring fairer evaluation metrics. Additionally, our
analysis of DPO demonstrates that its gains in human preference are influenced by response length,
underscoring the importance of unbiased evaluation benchmarks. Overall, AdapAlpaca provides a
robust tool for advancing reliable and equitable model evaluation.

486 REFERENCES 487

Josh Achiam. Ste	even Adler. Sandhi	ni Agarwa	ıl. Lam	a Ahmad. Ilge Ak	kava. Florencia	Leoni Ale-
man, Diogo Al	meida, Janko Alten	schmidt, S	Sam Al	tman, Shyamal Ana	adkat, et al. Gpt	t-4 technical
report. arXiv p	reprint arXiv:2303	.08774,20)23.			
AL@Mata Lla	ma 2 madel and	2024	LIDI	1a+++ + a = + / / a++ ha		- 11

491 492 493

488

489

490

Liama 3 model card. 2024. URL https://github.com/meta-llama/ AI@Meta. llama3/blob/main/MODEL_CARD.md.

- 494 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, 495 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi 496 Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng 497 Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi 498 Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang 499 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. arXiv preprint 500 arXiv:2309.16609, 2023. 501
- 502 Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. The pushshift reddit dataset. In Proceedings of the international AAAI conference on web and social media, volume 14, pp. 830-839, 2020. 504
- 505 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay 506 Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. Alpagasus: Training a better alpaca with 507 fewer data. arXiv preprint arXiv:2307.08701, 2023. 508
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, 509 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An 510 open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https: 511 //lmsys.org/blog/2023-03-30-vicuna/. 512
- 513 Qianlong Du, Chengqing Zong, and Jiajun Zhang. Mods: Model-oriented data selection for instruction tuning. arXiv preprint arXiv:2311.15653, 2023. 514
- 515 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled al-516 pacaeval: A simple way to debias automatic evaluators. arXiv preprint arXiv:2404.04475, 2024a. 517
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos 518 Guestrin, Percy S Liang, and Tatsunori B Hashimoto. Alpacafarm: A simulation framework for 519 methods that learn from human feedback. Advances in Neural Information Processing Systems, 520 36, 2024b. 521
- 522 Esin Durmus, Faisal Ladhak, and Tatsunori B. Hashimoto. Spurious correlations in reference-free 523 evaluation of text generation. In Annual Meeting of the Association for Computational Linguistics, 2022. URL https://api.semanticscholar.org/CorpusID:248300077. 524
- 525 Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large lan-526 guage models. In The Twelfth International Conference on Learning Representations, 2024. 527
- Zhengyu Hu, Jieyu Zhang, Zhihan Xiong, Alexander Ratner, Hui Xiong, and Ranjay Kr-528 ishna. Language model preference evaluation with multiple weak evaluators. arXiv preprint 529 arXiv:2410.12869, 2024. 530
- 531 Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep 532 Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. Camels in a changing climate: Enhancing Im adaptation with tulu 2. arXiv preprint arXiv:2311.10702, 2023.
- 534 Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop 535 Kang. Benchmarking cognitive biases in large language models as evaluators. arXiv preprint 536 arXiv:2309.17012, 2023.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, 538 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models for language modeling. arXiv preprint arXiv:2403.13787, 2024.

540 Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita 541 Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, et al. From generation to judgment: 542 Opportunities and challenges of llm-as-a-judge. arXiv preprint arXiv:2411.16594, 2024a. 543 Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 544 Camel: Communicative agents for "mind" exploration of large language model society. In Thirty-545 seventh Conference on Neural Information Processing Systems, 2023a. 546 547 Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun 548 Liu. Llms-as-judges: A comprehensive survey on llm-based evaluation methods. arXiv preprint 549 arXiv:2412.05579, 2024b. 550 Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi 551 Zhou, and Jing Xiao. From quantity to quality: Boosting llm performance with self-guided data 552 selection for instruction tuning. arXiv preprint arXiv:2308.12032, 2023b. 553 554 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy 555 Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 2023c. 556 Wei Liu, Yang Bai, Chengcheng Han, Rongxiang Weng, Jun Xu, Xuezhi Cao, Jingang Wang, 558 and Xunliang Cai. Length desensitization in directed preference optimization. arXiv preprint 559 arXiv:2409.06411, 2024. 560 Annie Louis and Ani Nenkova. Automatically assessing machine summary content without 561 a gold standard. Computational Linguistics, 39:267-300, 2013. URL https://api. 562 semanticscholar.org/CorpusID:17829732. 563 564 Swaroop Mishra, Anjana Arunkumar, Bhavdeep Sachdeva, Chris Bryan, and Chitta Baral. Dqi: 565 Measuring data quality in nlp. arXiv preprint arXiv:2005.00816, 2020. 566 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong 567 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-568 ton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan 569 Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback. 570 In NeurIPS, 2022. 571 572 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and 573 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. 574 In NeurIPS, 2023. 575 Abdul Sohail Rao, J. N. Kim, Meghana Kamineni, Minxia Pang, Wh Lie, and Marc D. Succi. 576 Evaluating chatgpt as an adjunct for radiologic decision-making. medRxiv : the preprint server 577 for health sciences, 2023. URL https://api.semanticscholar.org/CorpusID: 578 256626649. 579 Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing 580 Huang. Loose lips sink ships: Mitigating length bias in reinforcement learning from human 581 feedback. In Conference on Empirical Methods in Natural Language Processing, 2023. URL 582 https://api.semanticscholar.org/CorpusID:263831112. 583 584 Linxin Song, Jieyu Zhang, Lechao Cheng, Pengyuan Zhou, Tianyi Zhou, and Irene Li. Nlpbench: 585 Evaluating large language models on solving nlp problems. arXiv preprint arXiv:2309.15630, 586 2023. 587 John Von Neumann. Mathematische grundlagen der quantenmechanik. Springer-Verlag, 2013. 588 589 Thuy-Trang Vu, Xuanli He, Gholamreza Haffari, and Ehsan Shareghi. Koala: An index for quanti-590 fying overlaps with pre-training corpora. arXiv preprint arXiv:2303.14770, 2023. 591 Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and 592 Zhifang Sui. Large language models are not fair evaluators. arXiv preprint arXiv:2305.17926, 593 2023a.

- Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, et al. Pandalm: An automatic evaluation benchmark for llm instruction tuning optimization. arXiv preprint arXiv:2306.05087, 2023b.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. arXiv preprint arXiv:2212.10560, 2022.
- Lai Wei, Zhiquan Tan, Chenghai Li, Jindong Wang, and Weiran Huang. Large language model evaluation via matrix entropy. arXiv preprint arXiv:2401.17139, 2024.
 - Minghao Wu and Alham Fikri Aji. Style over substance: Evaluation biases for large language models. arXiv preprint arXiv:2307.03025, 2023.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. ArXiv, abs/2308.08155, 2023. URL https://api. semanticscholar.org/CorpusID:260925901.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. arXiv preprint arXiv:2304.12244, 2023.
- Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, et al. Justice or prejudice? quantifying biases in llm-as-a-judge. arXiv preprint arXiv:2410.02736, 2024.
- Qi Zhang, Yiming Zhang, Haobo Wang, and Junbo Zhao. Recost: External knowledge guided data-efficient instruction tuning. arXiv preprint arXiv:2402.17355, 2024.
- Hao Zhao, Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Long is more for alignment: A simple but tough-to-beat baseline for instruction fine-tuning. arXiv preprint arXiv:2402.04833, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haotong Zhang, Joseph Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. ArXiv, abs/2306.05685, 2023. URL https://api.semanticscholar.org/CorpusID:259129398.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. arXiv preprint arXiv:2305.11206, 2023.

IMPLEMENTATION DETAIL А 649

A.1 EXPERIMENT SETUP

652 In our experiments, we follow the setup in the AlpacaEval Leaderboard⁷, using the GPT-4 Preview 653 (11/06) as Baseline and the Annotator. The references to GPT-3.5, LLAMA3-70b, and Qwen1.5 654 72b in the main text denote gpt-3.5-turbo-0125, meta-llama/Meta-Llama-3-70B-Instruct⁸, and Qwen/Qwen1.5-72B-Chat⁹, respectively. Following previous work (Wei et al., 2024), we calcu-655 656 late conditional entropy using the method described in (Von Neumann, 2013).

A.2 DATASET

659 AlpacaEval (Dubois et al., 2024b) comprises 805 instructions, including 252 from the self-660 instruct test set (Wang et al., 2022), 188 from the Open Assistant (OASST) test set, 129 from An-661 thropic's helpful test set (Zhou et al., 2023), 80 from the Vicuna test set (Chiang et al., 2023), and 662 156 from the Koala test set (Vu et al., 2023).

663

648

650

651

657

658

664 (Zhou et al., 2023) compiles a training dataset of 1000 prompts and responses, designed LIMA 665 to ensure stylistic consistency in outputs while maintaining diverse inputs. It also provides an open-666 source test set of 300 prompts and a development set of 50. The dataset is sourced from a variety of 667 platforms, mainly community Q&A websites such as Stack Exchange, wikiHow, and the Pushshift Reddit Dataset (Baumgartner et al., 2020), along with manually curated examples. Within these 668 Q&A communities, highly upvoted answers on Reddit often have a humorous or trolling tone, re-669 quiring extra effort to align them with the intended helpful chat assistant style. In contrast, responses 670 from Stack Exchange and wikiHow naturally align with this style. The inclusion of human-authored 671 examples further enhances the dataset's diversity. Our research specifically utilizes the test set from 672 the LIMA dataset to evaluate our models.

673 674

680

Vicuna (Chiang et al., 2023) divides 80 test instructions into eight distinct categories: Fermi prob-675 lems, commonsense, roleplay scenarios, coding/math/writing tasks, counterfactuals, knowledge, and 676 generic questions. This categorization is intended to thoroughly evaluate multiple aspects of a chat-677 bot's performance. Prior research indicates that the Vicuna dataset generally includes instructions of 678 lower difficulty and complexity (Xu et al., 2023). In our study, we used the Vicuna test set to specif-679 ically evaluate the performance of large language models across these varied instruction categories.

Self-Instruct (Wang et al., 2022) consists of 252 human-created test instructions, each associated 681 with a carefully designed output. This test set is curated to reflect the real-world applicability of 682 instruction-following models, covering a broad spectrum of domains including email composition, 683 social media, productivity software, and coding. The test instructions vary in style and format, 684 incorporating different task lengths and diverse input/output types such as bullet lists, tables, code 685 snippets, and mathematical equations. We employed the Self-Instruct test set in our research to rigorously assess our model's capability to comply with precise instructions across these varied 687 domains.

688

689 Wizardlm (Xu et al., 2023) comprises a training set of 70k examples with varied complexities, 690 initiated from 52k instructional data provided by Alpaca. Following M = 4 evolutionary cycles, the 691 collection expands to 250k instructions. In each cycle, from the six newly generated prompts-five 692 via in-depth evolution and one through in-breadth evolution-one is chosen randomly for each instruction. ChatGPT then generates responses, resulting in $52 \times 4 \times 3 = 624$ k instruction-response 693 pairs. The training subset selected for the Evol-Instruct dataset contains 70k of these instructions. 694 The test set, which includes 218 instructions, is sourced from a variety of platforms such as open-695 source projects and online forums, encapsulating 29 unique skills identified from authentic human 696 tasks. These skills range from Coding Generation & Debugging to Reasoning, Mathematics, Writ-697 ing, Handling Complex Formats, and Mastery over Extensive Disciplines. In our study, we utilized 698 the Wizardlm test set to thoroughly evaluate our model's ability to adhere to detailed instructions. 699

⁷⁰⁰ 701

⁷https://tatsu-lab.github.io/alpaca_eval

⁸https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct

⁹https://huggingface.co/Qwen/Qwen1.5-72B-Chat

Table 6: Scores given by commonly used reward models to concise, detailed, and original responses from GPT-4. The analysis shows that the scores consistently decrease from detailed to concise responses, highlighting the length bias within the reward model.

	LLM Response					Rewa	ard Mode	el				Avg.
		Eurus	Grmdis	Grmsft	UniF	Debba	Bebla	FsfairRM	Gerew	Misrmr	InteRM	
	Concise	1.819	1.984	-2.919	0.064	2.229	4.159	-1.404	-0.456	5.661	0.426	1.156
	Origin	3.564	4.009	-0.505	2.901	3.305	5.142	1.830	1.066	9.440	1.558	3.231
	Detailed	3.986	4.646	1.458	3.263	3.759	5.450	2.684	2.630	10.616	2.416	4.090

(Vu et al., 2023) consists of 180 authentic user queries obtained from the Internet. These Koala queries cover a diverse array of topics and are generally characterized by a conversational tone, underscoring their applicability to real-world chat-based applications. To prevent test-set leakage, we exclude any query that achieves a BLEU score over 20% when compared to examples from our training set. Furthermore, we do not consider queries related to programming or non-English languages, as the capabilities of our crowd-sourced raters—who form our evaluation team—do not extend to effectively assessing such content. We have exclusively utilized the Koala test set to assess our model's capability to process and respond to genuine user inquiries in a conversational setting.

A.3 **INFORMATION ABOUT USE OF AI ASSISTANTS**

We use GPT-4 as an AI assistant during the preparation of this manuscript.

Dataset			Start
special-concise 🗸 🗸	parti		
Enable Translation			
Question			
What breed dog is smallest?			
		Right	
John f (b 3 Inden tel.		(t) kig and spically de nor weigh	mor that is pound (p. 7 kg) They are known for the s and a boy, kg even, and large reg. Chikubhas dhan bare bid personalities and can be guith field?
Left Win	Ur	ndo	Right Win
	Save P	rogress	
D Download export joon		E tabel	

Figure 5: Example of the evaluation interface used in the human study, showing two outputs for a single input query. Participants assessed which output more accurately addressed the query, demonstrating the interface's role in ensuring unbiased evaluation.

В HUMAN EVALUATION PROCESS

To ensure the robustness of our findings and complement the automated evaluations, a thorough human evaluation was conducted.

Participants. The human evaluation involved 25 participants, all of whom are professionals or researchers in the tech industry with specific expertise in language models. These individuals were carefully selected to represent a broad spectrum of perspectives and expertise levels, ranging from early-career to senior researchers. Each participant was assigned randomly to different segments of the dataset to ensure a balanced and unbiased input across all items evaluated.

761

Data Segmentation and Assignment. The dataset, comprising 805 responses generated for each prompt and compared against a default reference, was strategically divided into eight distinct parts, each containing approximately 100 responses. This division was structured to facilitate manageability and focus during the evaluation process. By dividing the dataset into smaller, more manageable segments, we aimed to optimize the evaluation process without overwhelming the evaluators, thus maintaining a high standard of analysis quality.

Each of these eight segments was then randomly assigned to five different participants. This approach ensured that every subset of the dataset was evaluated by multiple individuals, enhancing the reliability and diversity of perspectives in the assessment process. Random assignment of participants to each segment helped minimize any potential bias, providing a balanced evaluation across all parts of the dataset.

This method of segmenting the data and assigning evaluators ensured that each response received
sufficient attention, contributing to the robustness and credibility of the evaluation results. By implementing this straightforward and strategic approach to data handling and evaluator assignment, we maintained a high standard of reliability and fairness throughout the evaluation process.

777

Evaluation Interface. The evaluation was facilitated using a custom-built interface on Gradio ¹⁰,
an open platform known for its robustness in sharing interactive machine learning models. Detailed
instructions were provided to each participant to minimize user error and bias. The interface displayed questions along with two model outputs side-by-side, labeled "Left" and "Right," with their
positions randomized to prevent positional bias. Figure 5 illustrates this setup.

This comprehensive human evaluation process not only validated the effectiveness of our proposed
 methodologies but also provided critical insights that significantly enriched our understanding of
 automated metric evaluations.

786 787

788

C POTENTIAL NEGATIVE SOCIETAL IMPACTS

While this research contributes to reducing bias in language model evaluations, it is important to consider potential indirect societal impacts that might arise:

792 Dependence on Automated Decision-Making. This study's focus on enhancing the accuracy of
 793 automated evaluations may inadvertently promote an over-reliance on AI-driven decision-making
 794 processes. While beneficial in many respects, such reliance could diminish the value placed on
 795 human judgment and intuition in areas where nuanced understanding and ethical considerations are
 796 paramount.

797

Perception and Trust in AI. By highlighting the capabilities and improvements in AI evaluations,
there might be an overestimation of AI reliability and fairness among the public and policymakers.
This could lead to misplaced trust in AI systems, overlooking their limitations and the necessity for
continuous oversight and human intervention.

802 803

804

809

D LENGTH BIAS ORIGINATING FROM RLHF

We believe that the length bias observed in LLMs essentially originates from the RLHF (Ouyang et al., 2022) process. As shown in Figure 7, during the RLHF process, humans may generally prefer more detailed responses when labeling preference data. This leads to ranking data where longer responses are generally ranked higher than shorter ones, causing the reward model to learn this

¹⁰https://github.com/gradio-app/gradio



Figure 6: Analysis of the 14 commonly used preference datasets on Hugging Face. The analysis shows that the lengths of chosen responses are generally longer than those of rejected responses, indicating a length bias in human preference labeling.

spurious correlation and incorrectly assume that length is a factor in human preference. This bias is further propagated to the aligned model during the training process using the reward model.

To verify our idea, we first analyze 14 commonly used preference datasets in huggingface, shown in Figure 6. We found that the lengths of chosen responses are generally longer than those of rejected responses. As detailed in Table 6, we also analyze the scores given by 10 commonly used reward models (Lambert et al., 2024) to detailed, original, and concise responses from GPT-4. The detailed description of these three prompts can be found in Section 3.4. We find that the scores consistently decrease across all reward models. The details of these datasets and reward models can be found in Appendix E. However, attributing human preference solely to response length is an oversimplification, as length is merely a superficial factor in how humans judge the quality of a sentence.

E PREFERENCE DATASET AND REWARD MODELS

In this appendix, we provide detailed information about the preference datasets and reward models used in Appendix D.

E.1 PREFERENCE DATASETS

AnthropicHH ¹¹: The AnthropicHH dataset evaluates the ULMA technique by replacing positive samples in a preference dataset with high-quality 'golden' data from GPT-4, aiming to enhance alignment methods like RLHF, DPO, and ULMA.

¹¹https://huggingface.co/datasets/Unified-Language-Model-Alignment/ Anthropic_HH_Golden

866

867

868

870

871 872 873

874

879

880

881 882 883

884

885

886 887

888

889

890

891 892

893

894

895

896 897

899

900

901 902

907 908

909

910

911 912 913



Figure 7: RLHF process contributing to length bias in LLMs. Human labelers often prefer detailed responses, leading to ranking data where longer responses are ranked higher. This creates a spurious correlation that the reward model learns and propagates to the aligned model.

¹²: The HC3 dataset, presented in "How Close is ChatGPT to Human Experts? Comparison HC3 Corpus, Evaluation, and Detection," offers a pioneering human-ChatGPT comparison corpus. It enables nuanced evaluations of ChatGPT's performance and its closeness to human expert outputs.

SHP ¹³: The SHP dataset, from the Stanford Human Preferences project, collects 385K human preferences across 18 subject areas, utilizing naturally occurring human-written responses on Reddit to enhance RLHF reward models and NLG evaluation. This dataset emphasizes the utility of response helpfulness over harm reduction.

WebgptCom ¹⁴: The WebgptCom dataset comprises 19,578 comparisons from the WebGPT project, designed for reward modeling. It features pairs of model-generated answers to questions, each scored by humans to determine preference, supporting the training of a long-form question answering model aligned with human preferences.

Helpsteer ¹⁵: The Helpsteer dataset, utilized for refining reward models in conversational AI, in-898 cludes preference data distinguishing helpful from unhelpful responses. It consists of paired entries labeled as 'chosen' and 'rejected', with respective scores reflecting their utility. The dataset includes 37,131 examples in the training split, emphasizing its scale for robust model training.

¹⁶: The Preference700K dataset comprises 700,000 preference comparisons be-903 Preference700K tween two conversational responses, 'chosen' and 'rejected', related to the same prompt. This large-904 scale dataset is structured to train and evaluate models on their ability to discern more favorable 905 conversational outcomes based on user interaction dynamics. 906

¹⁷: The PreMix dataset features 528,029 comparisons from preprocessed preference PreMix datasets, focusing on dialogues structured with a 'chosen' and 'rejected' response based on the same prompt. This dataset aids in training models to discern the more favorable responses in conversational settings.

- ¹²https://huggingface.co/datasets/Hello-SimpleAI/HC3
- 914 ¹³https://huggingface.co/datasets/stanfordnlp/SHP
- ¹⁴https://huggingface.co/datasets/openai/webgpt_comparisons 915
- ¹⁵https://huggingface.co/datasets/RLHFlow/Helpsteer-preference-standard 916
- ¹⁶https://huggingface.co/datasets/hendrydong/preference_700K 917

¹⁷https://huggingface.co/datasets/weqweasdas/preference_dataset_mix2

Ultrafeedback ¹⁸: Ultrafeedback is an improved version of the original dataset, now cleaned and binarized using average preference ratings. It eliminates problematic data from earlier versions, notably those influenced by the TruthfulQA dataset, and removes contributions from ShareGPT sources, ensuring cleaner and more reliable data for fine-tuning conversational AI on preference discernment.

Argilla ¹⁹: The Argilla dataset is a refined version of the UltraFeedback dataset, used to train the Zephyr-7B- β model. This dataset features 64k prompts with binarized completions, categorizing the highest scored as 'chosen' and one of the remaining as 'rejected'. It supports various training techniques including supervised fine-tuning, preference modeling for reward systems, and generation techniques like rejection sampling.

Lmsys ²⁰: The Policy1.4b dataset incorporates labels from the AlpacaFarm dataset and utilizes generated answers from a 1.4 billion parameter Pythia policy model. Responses are evaluated using the 'reward-model-human' as a gold standard. This dataset is pivotal for refining AI policy models through precise human preference feedback.

934

952 953

954 955

956

957 958

923

929

Policy1.4b ²¹: The Prometheus2 dataset, transformed from the "prometheus-eval/Preference-Collection", is crafted to enhance fine-grained evaluation capabilities in language models. This dataset pairs instructions with two responses, scored and chosen based on preference, facilitating nuanced evaluation and comparison aligned with human judgment.

Prometheus2 ²²: The Prometheus2 dataset, transformed from the "prometheus-eval/Preference-Collection", is crafted to enhance fine-grained evaluation capabilities in language models. This dataset pairs instructions with two responses, scored and chosen based on preference, facilitating nuanced evaluation and comparison aligned with human judgment.

944
 945
 946
 946
 947
 947
 948
 949
 949
 949
 949
 949
 941
 941
 942
 942
 943
 944
 945
 945
 946
 947
 947
 946
 947
 947
 947
 947
 948
 949
 949
 949
 949
 949
 949
 941
 941
 942
 942
 943
 944
 945
 945
 946
 947
 947
 947
 947
 947
 947
 948
 949
 949
 949
 949
 949
 949
 949
 941
 941
 942
 942
 942
 944
 945
 945
 945
 946
 947
 947
 947
 947
 947
 947
 947
 947
 947
 947
 948
 949
 949
 949
 949
 949
 949
 941
 941
 942
 942
 942
 942
 943
 944
 944
 945
 945
 945
 946
 947
 947
 947
 947
 947
 947
 948
 948
 949
 949
 949
 949
 949
 949
 949
 949
 949
 949
 949
 949
 949
 949

948 ElixLatent ²⁴: The ElixLatent dataset, designed around GPT-4, serves as a resource for training and evaluating latent preference modeling. It provides pairs of latent responses ('yw' and 'yl') and their corresponding contexts ('x'), allowing researchers to explore the nuances of preference dynamics in generated text.

E.2 REWARD MODELS

Eurus ²⁵: Eurus is a reward model trained on UltraInteract, UltraFeedback, and UltraSafety datasets. It excels in complex reasoning tasks and outperforms larger models, including GPT-4, by significantly enhancing language models' reasoning capabilities.

Grmdis ²⁶: Generalizable Reward Model (GRM), uses hidden state regularization to enhance generalization in reward models for large language models (LLMs). Initially built on fixed weights from a Llama-3-based model and fine-tuned only on a reward head, it significantly improves on standard benchmarks, demonstrating enhanced reasoning and safety metrics over existing models.

```
970 <sup>24</sup>https://huggingface.co/datasets/Asap7772/elix_latent_p
971 <sup>25</sup>https://huggingface.co/openbmb/Eurus-RM-7b
```

```
<sup>26</sup>https://huggingface.co/Ray2333/GRM-llama3-8B-distill
```

972
 973
 973
 974
 975
 Grmsft ²⁷: It is part of the Generalizable Reward Model (GRM) series, aimed at enhancing LLMs through hidden state regularization. It excels across various complex evaluative tasks, outperforming other high-capacity models in reasoning and safety.

976 UniF ²⁸: It is a reward model finetuned on the 'llm-blender/Unified-Feedback' dataset using the
977 Mistral-7B-Instruct architecture. Achieving an accuracy of 0.7740 on test sets, it excels at modeling
978 human preferences. The model integrates diverse preference data from multiple sources, enhancing
979 its applicability in aligning LLMs to human judgments across various conversational contexts.

Debba ²⁹: Debba is a reward model utilizing Deberta-v3-base, trained to evaluate QA models and serve as a reward mechanism in RLHF by predicting which generated answer aligns better with human judgment. It is trained on datasets such as webgpt_comparisons, summarize_from_feedback, and synthetic-instruct-gptj-pairwise, ensuring a consistent validation approach across varying domains.

Bebla ³⁰: It is a reward model trained to assess the quality of responses in QA evaluations and to provide scoring in RLHF. It was developed with datasets such as webgpt_comparisons, summarize_from_feedback, and synthetic-instruct-gptj-pairwise, ensuring it can reliably predict human preferences across diverse contexts.

FsfairRM ³¹: It is designed for RLHF applications including PPO, iterative SFT, and iterative DPO. This state-of-the-art reward model is licensed under PKU-Alignment/PKU-SafeRLHF-30K, demonstrating high performance across diverse metrics like chat, safety, and reasoning in Reward-Bench.

996 32. Gerew It is trained using BT loss on the weqweas-997 das/preference_dataset_mixture2_and_safe_pku dataset. This model is designed for efficiently 998 evaluating and aligning LLMs, offering a baseline performance that is well-suited for smaller-scale applications requiring rapid assessment of language model outputs. 999

Misrmr ³³: It is a reward model tailored for iterative Synthetic Frontier Tuning (SFT) and Dynamic Policy Optimization (DPO). Trained to enhance language generation tasks, it supports fine grained reward modeling to improve the alignment and efficacy of language models in diverse applications.

InteRM ³⁴: It is a reward model trained on the foundation of InternLM2-Chat-1.8B-SFT. This model has been trained using over 2.4 million preference samples, both human-annotated and AI-synthesized, achieving outstanding performance while ensuring a balance between helpful and harmless.

1010

1011 F CASE STUDY

1012

980

To demonstrate the superiority of AdapAlpaca, we present a case study. In Figure 8, for the given instruction, we generate a redundant model answer (shown in the blue box). When evaluated using the current AlpacaEval response (shown in the red box), the annotator (i.e., GPT-4) selected this redundant answer, which is significantly unaligned from human preference, as the simplicity of the question does not warrant such extensive verbosity. The reason GPT-4 chose this answer is that the excessive length increases the information mass, artificially inflating the perceived quality.

²⁷https://huggingface.co/Ray2333/GRM-llama3-8B-sftreg

^{1020 &}lt;sup>28</sup>https://huggingface.co/Ray2333/reward-model-Mistral-7B-instruct-Unified-Feedback

^{1021 &}lt;sup>29</sup>https://huggingface.co/OpenAssistant/reward-model-deberta-v3-base

^{1022 &}lt;sup>30</sup>https://huggingface.co/OpenAssistant/reward-model-deberta-v3-large

^{1023 &}lt;sup>31</sup>https://huggingface.co/sfairXC/FsfairX-LLaMA3-RM-v0.1

^{1024 &}lt;sup>32</sup>https://huggingface.co/Ray2333/Gemma-2B-rewardmodel-baseline

^{1025 &}lt;sup>33</sup>https://huggingface.co/hendrydong/Mistral-RM-for-RAFT-GSHF-v0

³⁴https://huggingface.co/internlm/internlm2-1_8b-reward



Figure 8: Case study on comparing GPT-4 and human vote on AlpacaEval and AdapAlpaca. In
AlpacaEval, GPT-4 votes for the verbose answer, but humans vote for the concise reference answer,
while in AdapAlpaca, GPT-4 and humans vote for the same answer, demonstrating a better LLMhuman alignment on AdapAlpaca.

Table 7: The subscripts in the LCWR and WR columns indicate the differences between these metrics and the corresponding Human WR. A larger absolute value denotes a greater disparity between the annotator's evaluation and Human Preference. "LLM Response" denotes different responses to AlpacaEval questions for Llama3-70B, with detailed content available.

LLM Response		AlpacaEval	AdapAlpaca		
EEM Response	Human	LCWR	WR	Human	WR
Concise	5.67	$25.73_{\pm 20.06}$	$11.10_{\pm 5.43}$	7.24	$6.10_{-1.14}$
Detailed	46.12	$38.62_{-7.50}$	$50.80_{\pm 4.68}$	41.99	$42.98_{\pm 0.99}$
Quality Enhancement	53.48	$42.59_{-10.89}$	$56.50_{\pm 3.02}$	51.63	$50.89_{-0.74}$

1061

In contrast, when using AdapAlpaca, it allows us to control for content while varying the length, thereby isolating the effect of length from that of content quality.

G HUMAN STUDY WITH MORE MODEL

To provide a more comprehensive view of our human evaluation study, we conducted experiments on more LLMs, including Llama3-70B and Qwen1.5-72B. The results are summarized in Tables 7 and 8. These results further validate AdapAlpaca as a robust metric for aligning model evaluations with human preferences, effectively addressing the shortcomings of LCWR.

1067 H DATASET INFORMATION

1068

The data generation prompt, as outlined in Table 9, is carefully crafted to instruct GPT-4 to generate
 responses within predefined word limits. This prompt directed the model to generate content that is
 relevant to the given question and strictly adheres to the specified length constraints.

1072

1074

73 H.1 DATASET DOCUMENTATIONS.

The dataset comprises five JSON files for the *AdapAlpaca-200*, *AdapAlpaca-400*, *AdapAlpaca-600*, *AdapAlpaca-800*, and *AdapAlpaca-1000*. Each file is generated using our length control prompt technique with the Alpaca dataset employing the GPT-4 1106 model.

- 1078 Each data file contains a list of items with the following fields:
- 1079
- instruction: the prompt is given to generate the response.

1092

1093

1094 1095

1099 1100 1101

1102 1103

1104

1105

1106 1107

Table 8: The subscripts in the LCWR and WR columns indicate the differences between these metrics and the corresponding Human WR. A larger absolute value denotes a greater disparity between the annotator's evaluation and Human Preference. "LLM Response" denotes different responses to AlpacaEval questions for Qwen1.5-72B.

LIM Response		AlpacaEva	AdapAlpaca		
ELM Response	Human	LCWR	WR	Human	WR
Concise	9.24	$31.03_{+21.79}$	$13.20_{+3.96}$	8.56	$7.40_{-1.16}$
Detailed	45.52	$38.50_{-7.02}$	$42.70_{-2.82}$	39.97	$38.92_{-1.05}$
Quality Enhancement	46.61	$40.62_{-5.99}$	$48.87_{\pm 2.26}$	43.18	$44.01_{\pm 0.83}$

Table 9: Prompt for dataset generation, with $\{\max \text{ word}\}-\{\min \text{ word}\}\$ ranges set as 0-200, 200-400, 400-600, 600-800, and 800-1000.

Dataset generation prompt

You are a helpful assistant, highly attentive to the specified token range required from user. Respond to the following question, your reply must only be within $\{\max \text{ word}\}$ - $\{\min \text{ word}\}$ words.

- generator: identifies the model used.
- dataset: specifies the dataset used.
 - output_word_count: the word count of the generated response.
 - output: the actual text generated by the model.

1108 H.2 INTENDED USES.

The provided datasets, *AdapAlpaca-200*, *AdapAlpaca-400*, *AdapAlpaca-600*, *AdapAlpaca-800*, and *AdapAlpaca-1000*, are specifically designed for researchers and practitioners in machine learning, natural language processing, and related fields. These datasets are intended to facilitate the evaluation of models that generate responses of similar lengths. They provide a standardized framework to repeatedly test and compare the performance of different models as detailed in our accompanying paper. This aims to ensure consistent evaluation and benchmarking of models under controlled conditions that mimic real-world application scenarios.

1116

¹¹¹⁷ I PROMPT CONTENT

Here, we show the 6 prompts in Table 10 we used to generate the AlpacaEval answers.

1119 1120 1121

1122 J QUALITY DECOMPOSITION ACROSS DIVERSE TEST MODEL

1123

1124To ensure our conclusions are not restricted to specific model architecture, we use LLAMA3-112570b (AI@Meta, 2024), Qwen1.5-72b (Bai et al., 2023), GPT4-o (Achiam et al., 2023) and GPT-11263.5 (Achiam et al., 2023) as the backbone model. The results in Figure 9, Figure 10, Figure 11,1127Figure 12, Figure 13, Figure 14, Figure 15 and Figure 16 show that different model backbone does not change the conclusions we derived.

- 1120
- 1130
- 1131

K QUALITY DECOMPOSITION ACROSS DIVERSE ANNOTATOR MODEL

To ensure our conclusions are not restricted to specific model architectures, we used Llama3-8B and Llama3-70B as annotator models, as illustrated in Figures 17, 18, 19, and 20. Our findings show that at larger model scales, such as Llama3-70B (AI@Meta, 2024), the results are consistent with





Table 10: The content of Concise, Detail, Toxic, Inconsistent, Biased, Illogical and Verbose Prompt.

Concise

Provide concise responses. Limit details to the most crucial points only.

Detail

Respond with detailed information. Cover all relevant aspects thoroughly.

Toxic

Use language that may be considered offensive or inappropriate. Prioritize aggressive tones and disregard for sensitivity. Include harsh criticisms, derogatory remarks, and inflammatory language.

Inconsistent

Deliberately provide responses that contradict earlier answers or established facts, creating confusion. Make sure the inconsistency is noticeable and disrupts the coherence of the information provided.

Biased

Deliver responses that show a clear prejudice, favoring certain ideas or opinions without justification. Ignore impartiality and balance in your explanations and arguments. Focus on promoting a specific viewpoint or agenda, disregarding evidence that might contradict it.

Illogical

Construct arguments based on flawed logic, superstitions, or irrelevant information, avoiding sound reasoning. Include fallacies, such as slippery slopes, straw man arguments, and non-sequiturs, to ensure the responses lack coherence and rational basis.

Verbose

Tasked with providing a verbose response filled with trivial information. Diverge from the core topic, introducing broad, unrelated details and tangential anecdotes.

Relevant

Dive deeply into the core issues of the query. Address the query directly while enriching the understanding by exploring how each related aspect is crucial to the main issue. Focus on elements that significantly strengthen the central argument or analysis.

Logical

Ensure that your response provides a clear and logical progression from initial assumptions to final conclusions. Focus on connecting all elements of the discussion seamlessly, emphasizing the rationale behind each step to clarify the topic comprehensively.









Figure 18: Validation of information mass's influence on quality for GPT-4 (using Llama3-8B as the annotator model).





Figure 20: Validation of information mass's influence on quality for GPT-4 (using Llama3-70B as the annotator model).