Improving LLM-as-a-Judge Inference with the Judgment Distribution

Anonymous EMNLP submission

Abstract

Using language models to scalably approximate human preferences on text quality (LLMas-a-judge) has become a standard practice applicable to many tasks. A judgment is often extracted from the judge's textual output alone, typically with greedy decoding. However, LLM judges naturally provide *distributions* over judgment tokens, inviting a breadth of inference methods for extracting fine-grained preferences. We find that taking the mean of 011 the judgment distribution consistently outperforms taking the mode (i.e. greedy decoding) in all evaluation settings (i.e. pointwise, pairwise, and listwise). We further explore novel methods of deriving preferences from judgment distributions, and find that methods incorporat-018 ing risk aversion often improve performance. 019 Lastly, we analyze LLM-as-a-judge paired with chain-of-thought (CoT) prompting, showing that CoT can collapse the spread of the judgment distribution, often harming performance. Our findings show that leveraging distributional output improves LLM-as-a-judge, as opposed to using the text interface alone.

1 Introduction

027

034

042

LLM-as-a-judge has emerged as a scalable framework for evaluating model outputs by approximating human annotation (Lin et al., 2024; Li et al., 2024b; Dubois et al., 2024). Typically, such systems prompt off-the-shelf LLMs to score a response or rank multiple responses to a given user prompt. LLM-as-a-judge methods boast strong agreement with human judgments across a breadth of domains and criteria (Zheng et al., 2023b; Ye et al., 2023), despite current limitations (Koo et al., 2023; Tan et al., 2024).

Most prior work involving LLM-as-a-judge elicits judgments through the LLM's text interface (Lin et al., 2024; Zhu et al., 2023; Ye et al., 2023), where the most likely token (i.e. the mode of the next token distribution) or a sampled token is taken to represent the LLM's judgment. Recent works (Lee et al., 2024a; Liu et al., 2023b; Yasunaga et al., 2024) have suggested that taking the mean of the score token distribution can better represent the LLM's judgment. In this work, we comprehensively evaluate design choices for leveraging LLM judges' distributional output.¹ 043

045

047

050

051

052

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

078

079

We show that the mean consistently outperforms the mode in the pointwise, pairwise, and listwise settings (i.e. evaluating one, two, and many responses at a time). Specifically, the mean achieves higher accuracy in 42 out of 48 cases on RewardBench (Lambert et al., 2024) and MT-Bench (Zheng et al., 2023b). We further explore novel methods of deriving preferences from score distributions (Section 4). For example, incorporating risk aversion often improves performance. Categorizing methods as discrete or continuous, where discrete methods (e.g. mode) are simple to interpret like rubric scores, we find that continuous methods outperform discrete methods, due to the latter often predicting ties and failing to capture slight preferences. In particular, the mode assigns ties more frequently than every other method, leading to the lowest accuracy even among discrete methods.

We further study how chain-of-thought (CoT) prompting (Wei et al., 2022) impacts the performance of LLM-as-a-judge. After the CoT reasoning, LLMs often exhibit sharper score distributions, making the mean judgment similar to the mode. Removing CoT increases the spread of the judgment distribution, often improving performance, and more so for taking the mean than taking the mode (e.g. absolute +6.5% for mean vs. +1.4% for mode, on average with pointwise scoring on RewardBench), demonstrating the synergy between eliciting and using distributional output.

Our findings stress the importance of leveraging

¹We provide implementations of the evaluated methods at https://anonymous.4open.science/r/distributional-judge-D756/README.md.



Figure 1: Pointwise LLM judge's logits produce a score distribution. We show two ways to compare two score distributions: (1) comparing the modes of the distributions and (2) comparing the means of the distributions.

distributional output to maximize the effectiveness of LLM-as-a-judge, as opposed to using the text interface alone. As LLM-as-a-judge paradigms are widely adopted for complex tasks, improving best practices for using LLM-as-a-judge can impact many end tasks' development and evaluation.

Background 2

087

090

094

095

100

101

103

104

105

106

111

2.1 LLM-as-a-Judge Settings

We briefly review three settings for LLM-as-ajudge; see Appendix A for more background.

Pointwise Scoring The LLM judge scores the two texts independently on a scale from 1 to some K, as shown in Figure 1 (Zheng et al., 2023b; Lin et al., 2024; Cui et al., 2023).

Pairwise Scoring The LLM judge scores both texts in a single prompt (Zhu et al., 2023; Saha et al., 2023; Chan et al., 2023). To account for position bias, we prompt the LLM judge twice, once for each order of presentation, and average the outputs (Lee et al., 2024a).

Pairwise Ranking The LLM judge states which of the two texts it prefers (Lin et al., 2024; Li et al., 2024b; Dubois et al., 2024). As with pairwise scoring, we prompt the LLM judge twice, once for each order of presentation.

2.2 Related Work

Mean Judgment Several prior works have used 107 the mean of the judgment distribution, mostly in 108 the pointwise setting. Liu et al. (2023b); Lee et al. 109 110 (2024a); Saad-Falcon et al. (2024) note the benefits of the mean but do not empirically compare it with the mode. Zawistowski (2024), Hashemi 112 et al. (2024), Lukasik et al. (2024) show that the 113 mean outperforms the mode for summary scoring, 114

dialogue scoring, and other regression tasks. Concurrent work (Yasunaga et al., 2024) shows that the mean outperforms the mode on RewardBench (Lambert et al., 2024), but the paper's focus is on data-efficient alignment.

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

152

Lee et al. (2024a); Zhai et al. (2024) use pairwise judgment distributions to train a student model, but do not empirically compare with distillation using one-hot judgments. In this work, we benchmark the mode, the mean, and newly proposed methods for leveraging distributional judgments across the pointwise, pairwise, and listwise settings.

CoT Zheng et al. (2023b) presented preliminary evidence that CoT benefits LLM-as-a-judge. Other LLM-as-a-judge systems have been proposed that take advantage of LLMs' ability to perform CoT reasoning (Ankner et al., 2024; Feng et al., 2024). On the other hand, Liu et al. (2024f) evaluate many evaluation protocols and find that CoT can hurt performance. However, their analysis assumes access only to the judges' text interface, not examining the effect of CoT on the judgment distribution. In this work, we analyze the interplay between CoT and the inference method (e.g. mode vs. mean).

Related phenomena on the effect of CoT have been studied in the literature (Chiang and Lee, 2023; Stureborg et al., 2024; Liu et al., 2024a; Lee et al., 2023; Sprague et al., 2024; Hao et al., 2024; Zheng et al., 2023b). Wang and Zhou (2024) show the sharpening effect of CoT, which improves performance on numerical reasoning tasks. In this work, we show that this sharpening effect can be harmful when the LLM is used as a judge.

Distributional Reward Models Using distributional judgment makes it possible for LLM judges to represent pluralistically aligned preferences (Sorensen et al., 2024; Siththaranjan et al., 2023; Kumar et al., 2024). Compared to existing work on

distributional reward models (Siththaranjan et al., 2023; Zhang et al., 2024b; Li et al., 2024a; Dorka, 2024; Poddar et al., 2024; Padmakumar et al., 2024), (1) our setting involves LLMs not trained or prompted for distributional judgment (Meister et al., 2024), and (2) LLM judges can produce arbitrary distributions over a flexibly chosen discrete judgment space.

3 Distributional Judgment

In this section, we present our findings comparing mode vs. mean inference and CoT vs. no-CoT prompting for LLM-as-a-judge systems.

3.1 Methods

153

154

155

156

158

159

160

162

163

164

165

166

167 168

169

171

172

173

174

175

176

178

179

181

182

183

184

187

189

190

193

194

195

To infer a judgment from the LLM's output distribution, we use the mode or the mean. With **mode**, we perform greedy decoding to produce a judgment token and discard the logits. With **mean**, we compute a weighted average of the judgment options, weighting each judgment option by the probability assigned to its token. See Appendix B for details.

3.2 Experimental Setup

Models As LLM judges, we use gpt-4o-2024-08-06 (shortened to GPT-4o) (OpenAI et al., 2024), Llama-3.1-8B-Instruct (Llama-3.1-8B) (Dubey et al., 2024), Mistral-7B-Instruct-v0.3 (Mistral-7B) (Jiang et al., 2023), and Prometheus-2-7B (Kim et al., 2024). We cover a commonly used closedsource LLM² (GPT-4o), as well as smaller opensource variants.

Inference Settings We prompt the LLM judge with or without CoT reasoning, i.e. to provide a brief explanation before stating the judgment. We use greedy decoding for CoT prompting. See Appendix C for prompts.

We softmax the judgment logits into judgment probabilities with temperature 1. We use the score space $\{1, \ldots, K = 9\}$ in this section.

Evaluation Datasets and Metrics We evaluate on RewardBench (Lambert et al., 2024) and MT-Bench (Zheng et al., 2023b), two canonical datasets for preference modeling with human annotations. Each data instance contains a prompt, a preferred response, and a dispreferred response.

Model	Setting	Method	Reward Bench	MT-Bench
	point score	mode	85.1, 84.0	81.9, 80.5
Ģ	point score	mean	<u>87.4</u> , 88.0	83.6 , <u>83.2</u>
1-4	nair score	mode	86.7, <u>87.4</u>	<u>86.2, 86.5</u>
Ę.	pair score	mean	<u>87.1</u> , 87.6	<u>86.3</u> , 86.8
0	nair rank	mode	88.4, 89.7	86.3, 85.6
	pan rank	mean	88.6, 90.5	87.3 , 85.9
	noint soora	mode	69.6, 72.2	74.9, 71.9
- A	point score	mean	72.7, 79.3	78.7, 81.5
-8 II	noir cooro	mode	71.7, 75.2	82.6 , <u>82.4</u>
Lla 3.1	pair score	mean	72.1, 76.8	<u>82.3</u> , 81.2
· · I	pair rank	mode	68.9, 58.9	76.2, 63.0
		mean	74.2 , 68.6	80.0 , 76.5
		mode	60.4, 62.7	59.5, 66.2
7B	point score	mean	63.8, 72.1	62.6, 74.0
al-	nair score	mode	67.3, 68.9	<u>79.3, 79.8</u>
istr	pair score	mean	68.1, 71.0	<u>80.0</u> , 80.4
Σ	noir ronk	mode	56.3, 53.8	51.5, 51.5
	pail failk	mean	63.9 , 59.1	73.5 , 65.5
~	noint soora	mode	64.3, 66.0	72.5, 73.5
ena	point score	mean	64.6, 75.2	72.1, 81.6
7B	noir sooro	mode	71.0 , 68.7	78.4, <u>80.8</u>
¦⊃ŭ	pair score	mean	70.5, <u>70.8</u>	78.3, 80.9
Prc	pair rank	mode	59.6, 48.2	51.5, 43.0
	pan rank	mean	69.7 , 48.8	75.4 , 33.4

Table 1: Mode vs. mean and CoT vs. no-CoT (comma-separated) accuracy results (%). For each base model+setting, we bold the best result and underline results not significantly worse ($\alpha = 0.05$). The mean outperforms the mode in 42 out of 48 cases. No-CoT outperforms CoT in 14 out of 16 cases when using the mean for pointwise or pairwise scoring.

We evaluate accuracy on the binary classification task; predicting the correct winner, a tie, or the wrong winner gets 1, 0.5, or 0 points, respectively (Lambert et al., 2024). RewardBench contains 2,985 (prompt, response 1, response 2) triplets, each labeled with the preferred response. Since MT-Bench has multiple human judgments per triplet, we compute accuracy using only triplets with unanimous human judgments (1,132 out of 1,814). See Appendix D for dataset details.

3.3 Results

Table 1 shows our main results, comparing mode vs. mean and CoT vs. no-CoT across various prompt settings and LLMs.

Mean outperforms mode The mean outperforms the mode in 42 out of 48 cases. In Table 10, we provide a subset breakdown of RewardBench and observe particularly large gains for pointwise scoring on the Reasoning subset.

CoT often harms LLM-as-a-judge For the scoring settings, no-CoT outperforms CoT in 14 out of

²Many proprietary LLMs such as Google Gemini and Anthropic Claude do not provide logit access, preventing us from including them in our experiments. In Appendix E.1, we provide partial results for DeepSeek-V3, whose trends match those of GPT-40.

Model	Setting	RewardBench	MT-Bench
GPT-40	point score pair score pair rank	.039, .103 .042, .066 .002, .065	.041, .116 .038, .064 .012, .114
Llama -3.1-8B	point score pair score pair rank	.060, .101 .054, .106 .215, .318	.068, .093 .047, .092 .186, .331

Table 2: Average standard deviation of judgment distribution, with judgment options rescaled to [0, 1]. Comma-separated values in each cell are with and without CoT. No-CoT always has a greater standard deviation.

16 cases when using the mean. For the pairwise ranking setting, CoT outperforms no-CoT, except with GPT-40 on RewardBench.

217

219

221

222

226

227

We interpret the harmful effect of CoT on pointwise scoring with the smaller models as being due to *sharpening*, whereby the initial entropy in the judgment is lost as the model commits to one instantiation of a reasoning trace (Wang and Zhou, 2024). Table 2 confirms this trend by showing that the standard deviation of judgment distributions is lower for CoT than no-CoT. Moreover, removing CoT benefits the mean more than the mode (e.g. $69.6 \rightarrow 72.2$ for mode vs. $72.7 \rightarrow 79.3$ for mean, with Llama-3.1-8B on RewardBench), revealing the synergy between eliciting and utilizing distributional judgment.

Which setting works the best? Comparing different LLMs, we find GPT-40 performs better with pairwise judgment (e.g. 88.0 for pointwise scoring vs. 90.5 for pairwise ranking on RewardBench) as in prior work, but the smaller models often do 237 better with pointwise judgment and rely heavily on CoT for pairwise ranking (e.g. with Prometheus-2-239 7B on MT-Bench, 75.4 \rightarrow 33.4 when removing CoT 240 from pairwise ranking, compared to 81.6 with no-241 CoT pointwise scoring). We believe this is because 242 pairwise judgment demands a more powerful judge 243 to leverage the context. Thus, in pairwise ranking 244 with the smaller models, the reasoning gained by 245 CoT often outweighs the distributional signal lost in the process. Nonetheless, using pairwise scoring 247 (where assigning individual scores can be viewed as an intermediate reasoning step) rather than pairwise ranking can eliminate the need for CoT, and we recover much of the gap on RewardBench, and match or exceed pointwise performance on MT-Bench. 253

4 Study on Pointwise Scoring

Beyond the mode and mean discussed in prior work and the previous section, we further explore the design space of utilizing distributional output from LLM scorers. 254

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

282

283

285

286

287

290

291

292

293

294

295

296

297

298

299

300

301

302

Discrete vs. Continuous We say a method is *discrete* if it compares two score distributions by their independently assigned scores that take values in $\{1, \ldots, K\}$. Otherwise, we say it is *continuous*. Discrete scores are often desirable for interpretability (e.g. simple rubrics) but, by the pigeonhole principle, can often result in tied comparisons and fail to capture slight preferences.

Additional Metric: Mean Squared Error For our further analysis, we report mean squared error (MSE) in addition to accuracy. For target labels in $\{0, 1\}$ (a unanimously preferred response), MSE is equivalent to the Brier score. Accuracy incentivizes predicting a winner instead of a tie as long as oracle confidence is over 50%. In contrast, expected MSE is optimized by exactly predicting the oracle confidence, thus serving as a measure of a method's calibration given the judge's distributional output.

On MT-Bench, we generalize the label space to [0, 1] by averaging the human judgments, thus allowing us to evaluate MSE on the full dataset. In Appendix F.1, we analyze alignment between the judgment distributions of LLMs and those of humans (as opposed to the average or majority vote).

4.1 Methods

Table 3 lists our extended methods for comparing two score distributions. We motivate the newly introduced methods below and provide details in Appendix B.1.

Users often prefer discrete methods (e.g. mode) because they are simple to interpret, even if they have lower accuracy than continuous methods (e.g. mean). This motivates the question of where the mode (the status quo method) ranks among discrete methods. To answer this question, we compare the mode to other discrete methods: rounded mean, median, and first percentile (discussed in the next paragraph).

Humans exhibit risk aversion when making decisions. They often disprefer negative outcomes more strongly than they prefer positive outcomes. However, this disposition is not captured by the measures of central tendency discussed so far.

Name	Description	Definition of NAME $(X_1, X_2) \in [-1, 1]$ (higher says X_1 is better, lower says X_2 is better)	Discrete or Continuous
MODE	Mode	$\operatorname{sgn}(r_1 - r_2)$ with $r_i = \operatorname{arg} \max_k P(X_i = k)$	Discrete
MEAN	Mean	$\frac{\mathbb{E}(X_1 - X_2)}{\mathbb{E} X_1 - X_2 + \sigma(X_1 - X_2)}$	Continuous
[MEAN]	Rounded mean	$\operatorname{sgn}(r_1 - r_2)$ with $r_i = \operatorname{argmin}_k \mathbb{E}X_i - k $	Discrete
MEDI	Median	$sgn(r_1 - r_2)$ with $r_i = Q_{X_i}(0.5)$	Discrete
1 P	1st percentile	$sgn(r_1 - r_2)$ with $r_i = Q_{X_i}(0.01)$	Discrete
RAM	Risk-averse mean	$\operatorname{MEAN}(X_1 - \sigma(X_1), X_2 - \sigma(X_2))$	Continuous
QT PS	Quantiles Probability of superiority	$\int_0^1 \operatorname{sgn}(Q_{X_1}(p) - Q_{X_2}(p)) \mathrm{d}p P(X_1 > X_2) - P(X_1 < X_2)$	Continuous Continuous

Table 3: Methods of comparing two score distributions X_1, X_2 over K score options. sgn is the sign function. $Q_X(p)$ denotes the value at the *p*-quantile. $\sigma(X)$ denotes the standard deviation; $\sigma_-(X) = \sqrt{\mathbb{E}[\max(\mathbb{E}X - X, 0)^2]}$ denotes the lower semi-deviation, a risk measure (Bond and Satchell, 2002).

Model	Method	Rewar	dBench	MT-]	Bench
		$Acc \uparrow$	$MSE\downarrow$	$Acc \uparrow$	$MSE\downarrow$
	MODE	84.0	.118	80.5	.145
	MEAN	88.0	.102	<u>83.2</u>	.097
	[MEAN]	85.2	.109	80.2	.146
GPT 4o	MEDI	84.6	.112	80.2	.142
UF 1-40	1 P	84.3	.116	81.0	.138
	RAM	88.4	.100	83.4	.096
	QT	87.9	.096	<u>83.2</u>	.118
	PS	87.8	.096	<u>83.3</u>	.103
	MODE	72.2	.192	71.9	.142
	MEAN	79.3	.155	81.5	.104
	[MEAN]	75.0	.186	75.0	.145
Llama	MEDI	73.6	.191	73.9	.142
-3.1-8B	1 P	76.0	.183	79.2	.147
	RAM	79.9	.152	<u>81.4</u>	.102
	QT	79.0	.164	81.1	.116
	PS	78.9	.161	<u>81.4</u>	.110

Table 4: Pointwise results over methods. No-CoT (see Table 11 for CoT). Text styling follows Table 1.

Thus, we investigate whether incorporating the human disposition of risk aversion into LLM-as-ajudge inference methods improves alignment with human preferences. The methods 1P (discrete) and RAM (continuous) reflect risk aversion. 1P takes an approach contrary to MODE; instead of focusing on where the most mass lies, 1P assigns a low score if there is even a 1% chance of such a low score (Siththaranjan et al., 2023). RAM is MEAN but with each distribution shifted down by its risk σ_{-} .

We have used MODE to represent the status quo LLM-as-a-judge inference method, which uses greedy decoding to obtain a judgment token. However, some prior works use a positive temperature, e.g. to obtain varied CoT chains (Zhang et al., 2024a), in which case a sampled judgment token is decoded rather than the mode. To account for the random nature of sampling, we design the method PS as the difference in winrates over repeated pairs of samples from the LLM judge (Siththaranjan et al., 2023). QT generalizes MEDI and 1P by averaging the comparisons over all quantiles, and can be viewed as PS but with X_1 and X_2 positively monotonically correlated.

322

323

324

325

326

327

328

329

330

332

333

335

336

337

338

339

340

341

342

343

344

345

347

349

350

351

352

354

355

4.2 Results

Main Takeaways

- Table 4 shows that the top pointwise methods are the continuous ones (MEAN, RAM, QT, PS), in both accuracy and MSE, indicating that they should be chosen over discrete methods.
- Even among discrete methods, MODE has the lowest accuracy in 3 out of the 4 cases, indicating that the mode is a suboptimal choice even if discrete scores are desired.
- 1P often outperforms MEDI (e.g. 79.2 vs 73.9 accuracy with Llama-3.1-8B on MT-Bench), and RAM slightly outperforms MEAN (e.g. 79.9 vs. 79.3 accuracy with Llama-3.1-8B on Reward-Bench), suggesting that risk aversion can be help-ful for preference modeling.

Study: Score Granularity and Ties We show here that ties explain the finding above that the discrete methods fall behind the continuous ones, and we experiment with score granularity as a remedy.

Table 5 shows that the discrete methods predict ties on a significant number of instances, on which MEAN is still able to achieve nontrivial accuracy. On the other hand, we find that on instances where a discrete method does not predict a tie, it has similar accuracy to MEAN (not shown; see Table 12), indicating that the performance gap is well explained by ties. Nonetheless, tie behavior varies by method; MODE has the most ties and the highest MEAN ac-



Figure 2: Comparing pairwise LLM-as-a-judge prediction based on **when** to aggregate the two judgments, one from each response pair presentation order. Pre- vs. post-aggregation (bottom vs. top in figure) can be likened to mean vs. mode, as the former aggregates at the distribution level while the latter aggregates at the text level (if mode is used).

Model	Method	Tie	rate	MEAN's	accuracy
		K = 9	K = 99	K = 9	K = 99
	MODE	.17	.20	72	73
CDT 4a	[MEAN]	.16	.03	67	53
GP1-40	MEDI	.17	.09	70	62
	1 P	.16	.08	66	60
	MODE	.35	.24	69	70
Llama	[MEAN]	.26	.07	64	61
-3.1-8B	MEDI	.29	.11	67	67
	1 P	.23	.08	65	57

Table 5: Tie analysis for discrete pointwise methods on RewardBench using no-CoT (see Table 12 for CoT and Table 13 for MT-Bench). We report results with two score granularity levels (K). Tie rate is the proportion of instances where the method predicts a tie, over which we report MEAN's accuracy (%); excess of 50% or 75% indicates room for improving accuracy or MSE, respectively.

curacy, amounting to the most untapped signal for determining the better response.

357

362

363

375

Table 5 further shows that granularizing the score space from K = 9 to K = 99 improves the expressivity of the discrete methods (except for MODE), drastically reducing the rate of ties, while MEAN accuracies remain similar or decrease.

Table 6 expands on the comparison between K = 99 and K = 9, reporting results from the same setting in Table 4 except for the granularity scale. Consistent with our motivation, the discrete methods (except for MODE) improve in accuracy, rivaling the continuous methods. Although MODE somewhat makes up for its low accuracy with a lower MSE than most other discrete methods on MT-Bench, it suffers the highest MSE on Reward-Bench.

Taken together, Tables 4-6 show that even in use cases where discrete scores are desired, one should consider alternatives to the mode.

	Method	Rewar	dBench	MT-I	Bench
		Acc \uparrow	$MSE\downarrow$	Acc \uparrow	$MSE\downarrow$
	MODE	81.7-2.3	.134+.016	78.4-2.1	.158+.013
	MEAN	86.7 _{-1.3}	$.108_{+.006}$	82.9-0.3	$.099_{+.002}$
0	[MEAN]	$86.5_{+1.3}$	$.127_{+.018}$	$82.7_{+2.5}$	$.182_{+.036}$
Γ.4	MEDI	85.2+0.6	$.126_{+.014}$	81.5+1.3	$.170_{+.028}$
Ę.	1 P	$86.4_{+2.1}$	$.116_{+.000}$	$82.7_{+1.7}$.165+.027
0	RAM	86.7 _{-1.7}	.104 +.004	83.0 _{-0.4}	.098 +.002
	QT	<u>86.6</u> _{-1.3}	$.114_{+.018}$	$82.7_{-0.5}$	$.147_{+.029}$
	PS	<u>86.6</u> _{-1.2}	<u>.105</u> +.009	82.4-0.9	$.107_{+.004}$
	MODE	72.0-0.2	.221+.029	75.1+3.2	.169+.027
В	MEAN	$79.3_{+0.0}$	$.156_{+.001}$	$81.3_{-0.2}$	<u>.103</u> 001
-8	[MEAN]	$78.5_{+3.5}$	$.198_{+.012}$	$80.7_{+5.7}$	$.180_{+.035}$
ς.	MEDI	$76.5_{+2.9}$.207+.016	$80.1_{+6.2}$.161 _{+.019}
na-	1 P	$78.5_{+2.5}$	$.195_{+.012}$	$81.5_{+2.3}$	$.177_{+.030}$
lar	RAM	79.7 _{-0.2}	.152 _{+.000}	$81.1_{-0.3}$.102 _{+.000}
Ц	QT	$78.7_{-0.3}$	$.177_{+.013}$	81.3+0.2	$.143_{+.027}$
	PS	78.6-0.3	$.163_{+.002}$	81.8 _{+0.4}	$.111_{+.001}$

Table 6: Pointwise results over methods (K = 99). No-CoT (see Table 14 for CoT). Subscripts denote change from K = 9 (Table 4). Text styling follows Table 1.

Sensitivity to Score Granularity In Appendix F.2, we analyze the sensitivity of different methods to score granularity, and find theoretically and empirically that the mode is the most sensitive method.

376

377

378

381

382

384

385

386

387

390

5 Study on Pairwise Ranking

The judgment styles in Section 3's overview were scoring (Section 4) and ranking. In this section, we analyze design decisions for pairwise ranking, and in Section 6 listwise ranking.

5.1 Design Decisions

As we explain below, the pairwise ranking experiments in Table 1 used Likert-2, post-aggregation for the mode, and pre-aggregation for the mean. We now consider alternative choices (see Appendix B.2.2 for details).

462

463

Center	Agg.	Rewar	dBench	MT-Bench		
conter	Time	$Acc \uparrow$	$MSE\downarrow$	$Acc \uparrow$	$MSE\downarrow$	
mode	post	56.7	.240	57.5	.192	
	pre	73.1	.265	78.1	.222	
median	post	56.8	.240	57.5	.192	
	pre	72.9	.261	78.0	.218	
mean	post	<u>73.2</u>	.207	78.2	.144	
	pre	73.2	.222	78.1	.155	

Table 7: Pairwise ranking results over methods using Likert-3 comparing pre- and post-aggregation. All methods use Llama-3.1-8B, CoT (see Table 15 for GPT-40 and no-CoT). Text styling follows Table 1.

K	Reward	dBench	MT-Bench		
	$Acc \uparrow$	$MSE\downarrow$	$\operatorname{Acc}\uparrow$	$MSE\downarrow$	
2	74.2 , 68.6	.187 , .214	80.0 , 76.5	.126 , .135	
3	<u>73.2</u> , 66.3	.222, .240	78.1, 70.8	.155, .155	
5	70.0, 58.5	.215, .234	77.1, 64.8	.142, .153	

Table 8: Pairwise ranking results over Likert-*K* scales, using pre-aggregation mean. Llama-3.1-8B, CoT (see Table 16 for GPT-40 and no-CoT). Text styling follows Table 1.

Timing of aggregation and measure of central tendency Pairwise judgment suffers from position bias, i.e. the LLM judge's sensitivity to the order in which the evaluated texts are presented, which is usually addressed by prompting the LLM judge twice, once for each order of presentation (Lee et al., 2024a). We examine the remaining question of whether to aggregate the two judgments *before or after* computing the measure of central tendency (mode, median, or mean), as shown in Figure 2. Pre- vs. post-aggregation can be likened to mean vs. mode, as the former aggregates at the distribution level while the latter aggregates at the text level (if mode is used).

Granularity We prompt the judge to express its preference on a *K*-point Likert scale: [>, <] (Likert-2), [>, =, <] (Likert-3), or $[\gg, >, =, <, \ll]$ (Likert-5) (Liu et al., 2024b).

5.2 Methods Results

393

395

399

400

401

402

403

404

405

406

407

408

409

Table 7 shows that accuracy depends little on the 410 measure of central tendency and mostly on when 411 we aggregate, with aggregating first leading to 412 413 higher accuracy (as much as $56.7 \rightarrow 73.1$ using the mode on RewardBench). Considering that the tim-414 ing of aggregation does not affect accuracy if the 415 two runs agree, this shows that even for inconsis-416 tent judgments caused by position bias, there is 417

still valuable signal in the relative magnitudes of preference that we can leverage by aggregating first.

On the other hand, an intuitive explanation for why the measure of central tendency has little effect on accuracy is that the judgment space is small, so there is high correlation between the signs of the measures of central tendency. In fact, they are equivalent in the pre-aggregation Likert-2 setting.

Although aggregating first improves accuracy, it harms MSE for mode and median, which we attribute to the volatile prediction of a binary winner when faced with the uncertain situation of positional inconsistency. Nevertheless, the mean (with either pre- or post-aggregation) is among the top accuracy methods while outperforming all other methods on MSE. This demonstrates the calibration benefit of using the judgment distribution to produce a continuous prediction.

With GPT-40 (not shown; see Tables 15, 16), MSE is always minimized with no-CoT, highlighting the discord between CoT's sharpening effect and calibration. In Appendix F.3, we further analyze position bias and find that CoT increases the occurrence of severe position bias.

5.3 Granularity Results

Table 8 compares the Likert scales used in the pairwise ranking prompt. We find that **Likert-2 per-forms the best overall**, in line with the AlpacaEval methodology (Dubois et al., 2024) but deviating from WB-Reward and Arena-Hard-Auto (Lin et al., 2024; Li et al., 2024b), which use Likert-5.

6 Listwise Judgment

Listwise judgment is not as prevalent as pointwise or pairwise judgment, but it offers efficiency (Zhu et al., 2024) while granting the judge the maximal context for comparison (Buyl et al., 2023).

6.1 Judgment Spaces and Methods

We consider two prompts for eliciting listwise preferences over N texts (Appendix C.4). Prompt 1 is the one proposed by Zhu et al. (2024), which prompts to produce all $\binom{N}{2}$ pairwise preferences and then aggregate them into a sorted list. Prompt 2 skips the intermediate pairwise step and asks to directly produce the list (Liu et al., 2023a; Qin et al., 2023). We can then extract all pairwise³

³We retain the pairwise evaluation setup from previous sections; see Appendix D.1 for discussion.

464 preferences from one of the following judgment
465 spaces using the mode (textual output) or the mean
466 (distributional output).

467

468 469

470

471

472

473

474

475

476

477

478

479

480

481

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

504

505

506

507

508

510

511

- INTERM (Prompt 1): Intermediate pairwise preferences (Likert-3, no-CoT, only one of the two presentation orders), which we view as the reasoning process leading to the list. This efficiently extends pairwise ranking to the listwise setting, similar to batch prompting (Cheng et al., 2023).
 - LIST (Prompt 1): Final list. For MEAN, we use the probability distribution over text identifiers at each rank, inspired by Zhuang et al. (2023); Reddy et al. (2024). Specifically, at rank r, denote $p_r(i)$ as the probability of decoding text (identifier) *i*. Decoding text *i* at rank *r* implies that any text *j* not yet decoded will be decoded at a later rank and is thus worse than text *i*, and vice versa. Hence, we define MEAN $(i, j) \in [0, 1]$ as the average of $\frac{p_r(i)}{p_r(i)+p_r(j)}$ over the ranks *r* until *i* or *j* is decoded.
 - DIRECT LIST (Prompt 2): LIST but with Prompt 2 (no intermediate pairwise step).

6.2 Experimental Setup

Models Due to the context length required for listwise ranking and the difficulty of the task, we limit our evaluation to GPT-40. In preliminary experiments, we found poor performance with the smaller models, but in Appendix E.1 we show that DeepSeek-V3 exhibits similar trends to GPT-40.

Datasets We evaluate on Nectar (Zhu et al., 2024), RM-Bench (Liu et al., 2024c), and MT-Bench (Zheng et al., 2023b).

From Nectar, we use a random subset of 1,000 prompts, each with 7 responses. We discard the GPT-4 judgments included in the dataset and collect our own silver labels using GPT-40 with pairwise ranking (Likert-5, no-CoT, pre-aggregation, mean). RM-Bench contains 1,327 prompts, each with 3 chosen and 3 rejected responses, yielding 9 pairwise preference labels. MT-Bench contains 160 prompts, each with 6 responses. See Appendix D for dataset details.

6.3 Results

Table 9 compares mode and mean in the listwise judgment spaces. The two methods have similar accuracy, but the mean has much lower MSE.

We find DIRECT LIST to be the most accurate judgment space (notably, outperforming pointwise

Space	Method	Ne	Nectar R		RM-Bench		MT-Bench	
		Acc	MSE	Acc	MSE	Acc	MSE	
• ,	mode	80.4	.155	62.1	.339	80.8	.201	
merm	mean	80.4	.048	62.5	.243	<u>80.7</u>	.121	
list	mode	82.2	.156	62.4	.376	83.7	.189	
list	mean	<u>82.0</u>	.105	61.7	.317	<u>83.5</u>	.157	
direct list	mode	86.1	.138	69.9	.301	86.8	.168	
	mean	86.4	.087	69.4	.267	85.9	.133	

Table 9: Listwise results (GPT-40). Text styling follows Table 1.

scoring on MT-Bench; see Table 4), while INTERM has the lowest MSE. We hypothesize that DIRECT LIST outperforms LIST due to the intermediate pairwise comparisons playing a similar role to CoT in the pointwise and pairwise settings, where distributional output is captured most intactly without it. Even so, in Appendix F.3 we find DIRECT LIST to suffer the most position bias, consistent with Zhu et al. (2024), while INTERM has the least. 512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

7 Conclusion and Recommendations

We comprehensively evaluated design choices for leveraging LLM judges' distributional output. For pointwise scoring, we showed that continuous methods (e.g. mean) outperform discrete methods (especially the mode) due to ties. For pairwise ranking, we related the mean vs. mode comparison to pre- vs. post-aggregation of the two presentation orders' judgments. Although smaller LLM judges suffer heavily from inconsistent judgments due to position bias, pre-aggregation effectively leverages the relative magnitudes of preference.

We showed that CoT collapses the spread of the judgment distribution, often hurting performance. This applies even to the challenging setting of listwise ranking, where accuracy was maximized by directly predicting the list without an intermediate pairwise step. We hope that highlighting this limitation of CoT encourages the development of reasoning mechanisms that preserve output diversity and calibration for judgment and other subjective or open-ended tasks.

Recommendations We summarize our findings into guidelines for choosing judgment settings. Large judges like GPT-40 should use pairwise ranking no-CoT, or direct listwise ranking as an efficient alternative. Smaller judges like Llama-3.1-8B should use pointwise scoring no-CoT. The mean should be used instead of the mode, but these setting guidelines apply even if one uses the mode.

601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647

648

649

650

651

652

653

599

600

551 Limitations

560

561

563

577

581

582

583

584

590

591

592

593

594

595

597

598

552Downstream PerformanceIn this paper, we553evaluate LLM-as-a-judge design decisions by their554performance on preference modeling datasets.555However, this setup may not reveal downstream556impacts. We do not explore the impact of distribu-557tional judgments on reinforcement learning from558AI feedback (RLAIF) (Lee et al., 2024a) or human559decision making.

Training Our experiments involve off-the-shelf LLMs as judges without specific tuning. We do not explore training LLM judges to express distributional judgments (Saad-Falcon et al., 2024). Similarly, we exclude distributional reward models (Dorka, 2024) from the scope of our study.

566CoTWe conclude from our results that CoT of-567ten hurts judgment performance. However, we568only consider one prompt design per setting for569eliciting CoT reasoning (Appendix C) and do not570perform prompt optimization. Furthermore, we do571not consider more extensive test-time scaling, such572as asking the judge to produce its own reference573response (Zheng et al., 2023b) or aggregating many574CoT judgment runs (Zhang et al., 2024a; Stureborg575et al., 2024).

Natural Language Judgments A valuable aspect of LLM-as-a-judge is its ability to augment judgments with interpretable rationales (Mahan et al., 2024; Byun et al., 2024; Ye et al., 2024b; Cao et al., 2024). However, the distributional judgments we consider here are limited to those that are easily quantifiable, and we do not propose methods for leveraging distributional output over natural language feedback. While it is possible to continue decoding a rationale after the judgment, the rationale will be conditioned on the decoded judgment and not reflect the distribution over the unchosen judgment options. One approach could be to decode several rationales, each conditioned on a different judgment option.

References

- Zachary Ankner, Mansheej Paul, Brandon Cui, Jonathan D Chang, and Prithviraj Ammanabrolu. 2024. Critique-out-loud reward models. *arXiv* preprint arXiv:2408.11791.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova Dassarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al.

2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv*, abs/2204.05862.

- Shaun A. Bond and Stephen E. Satchell. 2002. Statistical properties of the sample semi-variance. *Applied Mathematical Finance*, 9:219 – 239.
- Maarten Buyl, Paul Missault, and Pierre-Antoine Sondag. 2023. Rankformer: Listwise learning-torank using listwide labels. *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- Ju-Seung Byun, Jiyun Chun, Jihyung Kil, and Andrew Perrault. 2024. Ares: Alternating reinforcement learning and supervised fine-tuning for enhanced multi-modal chain-of-thought reasoning through diverse ai feedback. *ArXiv*, abs/2407.00087.
- Maosong Cao, Alexander Lam, Haodong Duan, Hongwei Liu, Songyang Zhang, and Kai Chen. 2024. Compassjudger-1: All-in-one judge model helps model evaluation and evolution.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shan Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *ArXiv*, abs/2308.07201.
- Zhoujun Cheng, Jungo Kasai, and Tao Yu. 2023. Batch prompting: Efficient inference with large language model apis. In *Conference on Empirical Methods in Natural Language Processing*.
- Cheng-Han Chiang and Hunghuei Lee. 2023. A closer look into automatic evaluation using large language models. *ArXiv*, abs/2310.05657.
- Brian Conrey, James Gabbard, Katie Grant, Andrew Liu, and Kent E. Morrison. 2013. Intransitive dice. *Mathematics Magazine*, 89:133 143.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. *ArXiv*, abs/2310.01377.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, et al. 2025. Deepseek-v3 technical report. *Preprint*, arXiv:2412.19437.
- Nicolai Dorka. 2024. Quantile regression for distributional reward models in rlhf. *ArXiv*, abs/2409.10164.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.

- 669 670 671 676 677

- 697
- 700 701
- 703
- 706

- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. Kto: Model alignment as prospect theoretic optimization. arXiv preprint arXiv:2402.01306.
- Xidong Feng, Ziyu Wan, Mengyue Yang, Ziyan Wang, Girish A. Koushiks, Yali Du, Ying Wen, and Jun Wang. 2024. Natural language reinforcement learning. ArXiv, abs/2402.07157.
- Mark Finkelstein and Edward O. Thorp. 2006. Nontransitive dice with equal means.
- Nate Gillman, Daksh Aggarwal, Michael Freeman, Saurabh Singh, and Chen Sun. 2024. Fourier head: Helping large language models learn complex probability distributions.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training large language models to reason in a continuous latent space. Preprint, arXiv:2412.06769.
- Helia Hashemi, Jason Eisner, Corby Rosset, Benjamin Van Durme, and Chris Kedzie. 2024. Llmrubric: A multidimensional, calibrated approach to automated evaluation of natural language texts. In Annual Meeting of the Association for Computational Linguistics.
- Zhengyu Hu, Jieyu Zhang, Zhihan Xiong, Alexander J. Ratner, Hui Xiong, and Ranjay Krishna. 2024. Language model preference evaluation with multiple weak evaluators. ArXiv, abs/2410.12869.
- Hawon Jeong, ChaeHun Park, Jimin Hong, and Jaegul Choo. 2024. Prepair: Pointwise reasoning enhance pairwise evaluating for robust instruction-following assessments. arXiv preprint arXiv:2406.12319.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. Preprint, arXiv:2310.06825.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. arXiv preprint arXiv:2405.01535.
- Alexander Y Klimenko. 2015. Intransitivity in theory and in the real world. Entropy, 17(6):4364-4412.
- Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2023. Benchmarking cognitive biases in large language models as evaluators. arXiv preprint arXiv:2309.17012.
- Sachin Kumar, Chan Young Park, Yulia Tsvetkov, Noah A. Smith, and Hanna Hajishirzi. 2024. Compo: Community preferences for language model personalization.

Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. arXiv preprint arXiv:2403.13787.

707

708

713

715

716

717

718

720

721

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

- Gyeong-Geon Lee, Ehsan Latif, Xuansheng Wu, Ninghao Liu, and Xiaoming Zhai. 2023. Applying large language models and chain-of-thought for automatic scoring. ArXiv, abs/2312.03748.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. 2024a. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. In Forty-first International Conference on Machine Learning.
- Noah Lee, Jiwoo Hong, and James Thorne. 2024b. Evaluating the consistency of llm evaluators.
- Dexun Li, Cong Zhang, Kuicai Dong, Derrick-Goh-Xin Deik, Ruiming Tang, and Yong Liu. 2024a. Aligning crowd feedback via distributional preference reward modeling. ArXiv, abs/2402.09764.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. 2024b. From crowdsourced data to highquality benchmarks: Arena-hard and benchbuilder pipeline. arXiv preprint arXiv:2406.11939.
- Weixian Waylon Li, Yftah Ziser, Yifei Xie, Shay B. Cohen, and Tiejun Ma. 2024c. Tsprank: Bridging pairwise and listwise methods with a bilinear travelling salesman model.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let's verify step by step. ArXiv, abs/2305.20050.
- Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. 2024. Wildbench: Benchmarking llms with challenging tasks from real users in the wild. arXiv preprint arXiv:2406.04770.
- Minqian Liu, Ying Shen, Zhiyang Xu, Yixin Cao, Eunah Cho, Vaibhav Kumar, Reza Ghanadan, and Lifu Huang. 2023a. X-eval: Generalizable multi-aspect text evaluation via augmented instruction tuning with auxiliary evaluation aspects. In North American Chapter of the Association for Computational Linguistics.

- 761 762 765 770 771 772 773 774 776 780 790 791 794 798 799 810 811 812
- 813
- 814 815

- Ryan Liu, Jiayi Geng, Addison J. Wu, Ilia Sucholutsky, Tania Lombrozo, and Thomas L. Griffiths. 2024a. Mind your step (by step): Chain-of-thought can reduce performance on tasks where thinking makes humans worse.
- Shang Liu, Yu Pan, Guanting Chen, and Xiaocheng Li. 2024b. Reward modeling with ordinal feedback: Wisdom of the crowd.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. G-eval: Nlg evaluation using gpt-4 with better human alignment. arXiv preprint arXiv:2303.16634.
- Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. 2024c. Rm-bench: Benchmarking reward models of language models with subtlety and style.
- Yinhong Liu, Zhijiang Guo, Tianya Liang, Ehsan Shareghi, Ivan Vuli'c, and Nigel Collier. 2024d. Aligning with logic: Measuring, evaluating and improving logical consistency in large language models. ArXiv, abs/2410.02205.
- Yinhong Liu, Han Zhou, Zhijiang Guo, Ehsan Shareghi, Ivan Vulic, Anna Korhonen, and Nigel Collier. 2024e. Aligning with human judgement: The role of pairwise preference in large language model evaluators. arXiv preprint arXiv:2403.16950.
- Yixin Liu, Kejian Shi, Alexander Fabbri, Yilun Zhao, Peifeng Wang, Chien-Sheng Wu, Shafiq Joty, and Arman Cohan. 2024f. Reife: Re-evaluating instructionfollowing evaluation. ArXiv, abs/2410.07069.
- Adian Liusie, Potsawee Manakul, and Mark JF Gales. 2023. Zero-shot nlg evaluation through pairware comparisons with llms. arXiv preprint arXiv:2307.07889.
- Adian Liusie, Vatsal Raina, Yassir Fathullah, and Mark J. F. Gales. 2024. Efficient llm comparative assessment: A product of experts framework for pairwise comparisons. ArXiv, abs/2405.05894.
- Charles Lovering, Michael Krumdick, Viet Dac Lai, Nilesh Kumar, Varshini Reddy, Rik Koncel-Kedziorski, and Chris Tanner. 2024. Are language model logits calibrated? ArXiv, abs/2410.16007.
- Michal Lukasik, Harikrishna Narasimhan, Aditya Krishna Menon, Felix X. Yu, and Sanjiv Kumar. 2024. Regression aware inference with llms. In Conference on Empirical Methods in Natural Language Processing.
- Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato, Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. 2024. Generative reward models. arXiv preprint arXiv:2410.12832.
- Nicole Meister, Carlos Guestrin, and Tatsunori Hashimoto. 2024. Benchmarking distributional alignment of large language models.

Niklas Muennighoff, Qian Liu, Qi Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro von Werra, et al. 2023. Octopack: Instruction tuning code large language models. ArXiv, abs/2308.07124.

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, et al. 2024. Gpt-4o system card. Preprint, arXiv:2410.21276.
- Vishakh Padmakumar, Chuanyang Jin, Hannah Rose Kirk, and He He. 2024. Beyond the binary: Capturing diverse preferences with reward regularization. ArXiv, abs/2412.03822.
- John W. Payne. 1976. Task complexity and contingent processing in decision making: An information search and protocol analysis. Organizational Behavior and Human Performance, 16(2):366-387.
- Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. 2024. Personalizing reinforcement learning from human feedback with variational preference learning. ArXiv, abs/2408.10075.
- Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, et al. 2023. Large language models are effective text rankers with pairwise ranking prompting. arXiv preprint arXiv:2306.17563.
- Zhen Qin, Junru Wu, Jiaming Shen, Tianqi Liu, and Xuanhui Wang. 2024. Lampo: Large language models as preference machines for few-shot ordinal classification. ArXiv, abs/2408.03359.
- Vyas Raina, Adian Liusie, and Mark Gales. 2024. Is llm-as-a-judge robust? investigating universal adversarial attacks on zero-shot llm assessment. arXiv preprint arXiv:2402.14016.
- Revanth Gangi Reddy, JaeHyeok Doo, Yifei Xu, Md Arafat Sultan, Deevya Swain, Avirup Sil, and Heng Ji. 2024. First: Faster improved listwise reranking with single token decoding. arXiv preprint arXiv:2406.15657.
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2023. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. ArXiv, abs/2308.01263.
- Jon Saad-Falcon, Rajan Vivek, William Berrios, Nandita Shankar Naik, Matija Franklin, Bertie Vidgen, Amanpreet Singh, Douwe Kiela, and Shikib Mehri. 2024. Lmunit: Fine-grained evaluation with natural language unit tests. ArXiv, abs/2412.13091.
- Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li. 2023. Branchsolve-merge improves large language model evaluation and generation. ArXiv, abs/2310.15123.

976

977

978

979

925

Richard P. Savage. 1994. The paradox of nontransitive dice. *American Mathematical Monthly*, 101:429–436.

870

871

874

875

878

879

882

883

898

899

900

901

902

903

904

905

906

907

908

909

910

911 912

913

915

916

917

918

919

921

922

923

- Markus Schulze. 2011. A new monotonic, cloneindependent, reversal symmetric, and condorcetconsistent single-winner election method. *Social Choice and Welfare*, 36(2):267–303.
- Ammar Shaikh, Raj Abhijit Dandekar, Sreedath Panat, and Raj Abhijit Dandekar. 2024. Cbeval: A framework for evaluating and interpreting cognitive biases in llms. *ArXiv*, abs/2412.03605.
 - Lin Shi, Weicheng Ma, and Soroush Vosoughi. 2024. Judging the judges: A systematic investigation of position bias in pairwise comparative assessments by llms. *arXiv preprint arXiv:2406.07791*.
 - Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. 2023. Distributional preference learning: Understanding and accounting for hidden context in rlhf. *arXiv preprint arXiv:2312.08358*.
 - Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell L Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, et al. 2024. Position: A roadmap to pluralistic alignment. In *Forty-first International Conference on Machine Learning*.
 - Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2024. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. *ArXiv*, abs/2409.12183.
 - Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. Large language models are inconsistent and biased evaluators. *ArXiv*, abs/2405.01724.
 - Sijun Tan, Siyuan Zhuang, Kyle Montgomery, William Y Tang, Alejandro Cuadron, Chenguang Wang, Raluca Ada Popa, and Ion Stoica. 2024. Judgebench: A benchmark for evaluating llm-based judges. *arXiv preprint arXiv:2410.12784*.
 - Raphael Tang, Xinyu Crystina Zhang, Xueguang Ma, Jimmy Lin, and Ferhan Ture. 2023. Found in the middle: Permutation self-consistency improves listwise ranking in large language models. In North American Chapter of the Association for Computational Linguistics.
 - T. N. Tideman. 1987. Independence of clones as a criterion for voting rules. *Social Choice and Welfare*, 4(3):185–206.
- Xuezhi Wang and Denny Zhou. 2024. Chain-ofthought reasoning without prompting. *ArXiv*, abs/2402.10200.
 - Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. 2023. Do-not-answer: A dataset for evaluating safeguards in llms. *ArXiv*, abs/2308.13387.

- Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii Kuchaiev, and Yi Dong. 2024a. Helpsteer2preference: Complementing ratings with preferences. *ArXiv*, abs/2410.01257.
- Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024b. Helpsteer2: Open-source dataset for training top-performing reward models. *ArXiv*, abs/2406.08673.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903.
- Michihiro Yasunaga, Leonid Shamis, Chunting Zhou, Andrew Cohen, Jason Weston, Luke Zettlemoyer, and Marjan Ghazvininejad. 2024. Alma: Alignment with minimal annotation.
- Chen Ye, Wei Xiong, Yuheng Zhang, Nan Jiang, and Tong Zhang. 2024a. Online iterative reinforcement learning from human feedback with general preference model.
- Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. 2023. Flask: Fine-grained language model evaluation based on alignment skill sets. *arXiv preprint arXiv:2307.10928*.
- Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yujia Zhou, Wei Shen, Dong Yan, and Yiqun Liu. 2024b. Beyond scalar reward model: Learning generative judge from preference data. *ArXiv*, abs/2410.03742.
- Krystian Zawistowski. 2024. Unused information in token probability distribution of generative llm: improving llm reading comprehension through calculation of expected values. *arXiv preprint arXiv*:2406.10267.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2023. Evaluating large language models at evaluating instruction following. *arXiv preprint arXiv:2310.07641*.
- Yuanzhao Zhai, Zhuo Zhang, Kele Xu, Hanyang Peng, Yue Yu, Dawei Feng, Cheng Yang, Bo Ding, and Huaimin Wang. 2024. Online self-preferring language models. *ArXiv*, abs/2405.14103.
- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2024a. Generative verifiers: Reward modeling as next-token prediction. *arXiv preprint arXiv:2408.15240*.
- Michael JQ Zhang, Zhilin Wang, Jena D Hwang, Yi Dong, Olivier Delalleau, Yejin Choi, Eunsol Choi, Xiang Ren, and Valentina Pyatkin. 2024b. Diverging preferences: When do annotators disagree and do models know? *arXiv preprint arXiv:2410.14632*.

981

983

985

986

987

988

993

- 999 1000 1001
- 1002 1003
- 1004 1005 1006
- 1007

1016

1019

1021

1022

1024

1025

1026

1028

1029

1030

1032

1017

1015

1014

1012 1013

Xu Zhang, Xunjian Yin, and Xiaojun Wan. 2024c. Contrasolver: Self-alignment of language models by resolving internal preference contradictions. ArXiv, abs/2406.08842.

- Yifan Zhang, Ge Zhang, Yue Wu, Kangping Xu, and Quanquan Gu. 2024d. General preference modeling with preference representations for aligning language models. ArXiv, abs/2410.02197.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric P. Xing, et al. 2023a. Lmsyschat-1m: A large-scale real-world llm conversation dataset. ArXiv, abs/2309.11998.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023b. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems, 36:46595-46623.
- Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, Karthik Ganesan, Wei-Lin Chiang, Jian Zhang, and Jiantao Jiao. 2024. Starling-7b: Improving helpfulness and harmlessness with rlaif. In First Conference on Language Modeling.
- Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2023. Judgelm: Fine-tuned large language models are scalable judges. arXiv preprint arXiv:2310.17631.
- Shengyao Zhuang, Honglei Zhuang, Bevan Koopman, and G. Zuccon. 2023. A setwise approach for effective and highly efficient zero-shot ranking with large language models. In Annual International ACM SI-GIR Conference on Research and Development in Information Retrieval.

Related Work: LLM-as-a-judge A Settings

LLM-as-a-judge has been used in pointwise (evaluating one response at a time), pairwise (two), and listwise (many) settings.

Pairwise judgment has the advantage of grounding each evaluated response in the other, creating for a more calibrated task and leading to better agreement with humans (Liusie et al., 2023). However, due to intransitivity in pairwise preferences (Liu et al., 2024e), the cost to sort N texts is $O(N^2)$ rather than $O(N \log N)$, compared to O(N) in the pointwise setting. In addition, pairwise comparisons are susceptible to position bias (Shi et al., 2024), which often must be addressed by running both orders and aggregating the results (Zeng et al., 2023; Li et al., 2024b). Pairwise comparisons have also been shown to be more biased toward superficial traits such as verbosity and tone, in both LLM

and human judges (Payne, 1976; Jeong et al., 2024), although pointwise scoring more easily falls victim to adversarial responses (Raina et al., 2024).

The listwise setting provides the maximal amount of context to the judge while keeping the same compute complexity as the pointwise setting. However, the judgment task becomes much more challenging (Qin et al., 2023; Koo et al., 2023), especially due to the amplified position bias (Zhu et al., 2024), and the combinatorially many orders makes it severely more daunting to address than in the pairwise case (Tang et al., 2023; Qin et al., 2024). To mitigate position bias, Zhu et al. (2024) leverage intermediate pairwise preferences for aggregation into a sorted list. Zhuang et al. (2023); Reddy et al. (2024) use the distribution from a single output token for listwise passage reranking, a related task to LLM-as-a-judge.

B Methods

Let A_1 and A_2 be two texts to compare. We describe the methods of predicting a value in [-1, 1]that signifies the advantage of A_1 over A_2 . For accuracy, we take the sign of the prediction. For MSE, we rescale predictions from [-1, 1] to [0, 1].

The prompts for the various settings are in Appendix C.

B.1 Pointwise Methods

We elaborate on the pointwise methods introduced in Section 4.1. The LLM judge independently judges A_1 and A_2 , producing score distributions over $\{1, \ldots, K\}$ for an integer K that define independent random variables X_1 and X_2 , which are used to compare A_1 and A_2 .

The methods are invariant to scaling and translating the judgment space, and all methods that do not take expectations \mathbb{E} (which assumes linearity) are invariant to applying a positive monotone transformation to the judgment space. The methods are all equivalent if the distributions are deterministic, thus our experiments evaluate their ability to leverage the LLM judge's distributional output.

The denominator in MEAN normalizes it into [-1,1], similar to sgn $(x) = \frac{x}{|x|}$, taking $\frac{0}{0}$ to be 0. The σ term lowers the magnitude of the prediction in the presence of uncertainty in a continuous manner. Specifically, let $k, k' \in \{1, \ldots, K\}$ with $k \neq k'$. For $\epsilon \in [0,1]$, let X_1 have a two-point distribution $(1 - \epsilon)\delta_k + \epsilon \delta_{k'}$ and let X_2 have a deterministic distribution δ_k . Then MEAN (X_1, X_2)

1035 1036 1037

1038

1039

1033

1034

1045 1046

1048

1050 1051

1052

1053

1054 1055 1056

1057 1058

1059

1060

1061

1062

1063

1064

1065

1066

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

as a function of ϵ is continuous at $\epsilon = 0$.

For MEAN, RAM, and PS, we assume X_1 and X_2 to be independent, but QT can be viewed as PS but with X_1 and X_2 positively monotonically correlated. By incorporating the sign function, QT and PS are less sensitive to extremal values than MEAN. In addition, QT and PS can model intransitive preferences, e.g. $PS(X_1, X_2)$, $PS(X_2, X_3) > 0 \Rightarrow PS(X_1, X_3) > 0$, which we analyze in Appendix F.4.

B.2 Pairwise Methods

1082

1083

1084

1085 1086

1087

1091

1092 1093

1094

1095

1096

1097

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1128

In the pairwise setting, we consider two prompting approaches for jointly evaluating the two texts A_1 and A_2 : scoring both texts (§B.2.1) and expressing a preference (§B.2.2).

To account for position bias, we prompt the LLM judge once for each order of presentation. For an order $\mathbf{o} \in \mathbf{O} := \{(1, 2), (2, 1)\}$, we use \mathbf{o} to denote dependence on the order (A_{o_1}, A_{o_2}) in which the texts appear in the prompt.

B.2.1 Pairwise Scoring

For a given order o, the LLM judge scores the two texts jointly in the same run. If we could obtain the joint distribution $P(X_{o_1}^{o}, X_{o_2}^{o})$, we could compute the marginals and use any method in Table 3. However, the judge first outputs the score for A_{o_1} and conditions on it when outputting the score for A_{o_2} , i.e. $X_{o_1}^{o}$ and $X_{o_2}^{o}$ are not independent. Thus, the full joint distribution $P(x_{o_1}, x_{o_2}) =$ $P(x_{o_1})P(x_{o_2} \mid x_{o_1})$ can only be obtained by injecting each $x_{o_1} \in \{1, \ldots, K\}$ into the context to access $P(x_{o_2} \mid x_{o_1})$. This is feasible with local models but not with API-access models where inference cost scales with K. Hence, we stick to a single run and condition on the greedily decoded $x_{o_1} = \arg \max_k P(X_{o_1}^{o} = k)$, giving us

1118
$$X_{\Delta}^{\mathbf{o}} \stackrel{d}{=} (X_1^{\mathbf{o}} - X_2^{\mathbf{o}}) \mid (X_{o_1}^{\mathbf{o}} = x_{o_1})$$

as a proxy for the score difference $X_1^{\mathbf{o}} - X_2^{\mathbf{o}}$. Se-1119 mantically, $X^{\mathbf{o}}_{\Delta}$ is symmetric (i.e. there should be 1120 no prior preference for A_1 or A_2), so we would like 1121 our scalar judgment to be some measure of *central* 1122 tendency (mode, median, or mean). As shown in 1123 Figure 2, we also have the choice of whether to 1124 aggregate the judgments from the two orders of 1125 presentation before or after computing the measure 1126 of central tendency. 1127

For pre-aggregation, we simply take the mixture

distribution,

$$P(X_{\Delta} = \delta) \coloneqq \frac{1}{|\mathbf{O}|} \sum_{\mathbf{o} \in \mathbf{O}} P(X_{\Delta}^{\mathbf{o}} = \delta)$$
 1130

for all $\delta \in \{-(K-1), \dots, K-1\}$, leaving more 1131 sophisticated approaches such as the convolution 1132 and Wasserstein barycenter for future study: 1133

$$\mathsf{AGG}\operatorname{-}\mathsf{MODE} \coloneqq \operatorname{sgn}(\operatorname{mode}(X_\Delta))$$
 1134

$$AGG-MEDI \coloneqq \operatorname{sgn}(\operatorname{median}(X_{\Delta}))$$
 1135

$$\operatorname{AGG-MEAN} \coloneqq \operatorname{MEAN}(X_\Delta),$$
 1136

where MEAN is defined as in Table 3, overloaded to take a single argument representing $X_1 - X_2$.

For post-aggregation, we sum the two scalar judgments from the two orders and normalize:

$$\text{MODE-AGG} \coloneqq \frac{\sum_{\mathbf{o} \in \mathbf{O}} \text{mode}(X_{\Delta}^{\mathbf{o}})}{\sum_{\mathbf{o} \in \mathbf{O}} |\text{mode}(X_{\Delta}^{\mathbf{o}})|}$$
114

$$\text{MEDI-AGG} \coloneqq \frac{\sum_{\mathbf{o} \in \mathbf{O}} \operatorname{median}(X_{\Delta}^{\mathbf{o}})}{\sum_{\mathbf{o} \in \mathbf{O}} |\operatorname{median}(X_{\Delta}^{\mathbf{o}})|}$$
 1142

$$\text{MEAN-AGG} := \frac{1}{|\mathbf{O}|} \sum_{\mathbf{o} \in \mathbf{O}} \text{MEAN}(X_{\Delta}^{\mathbf{o}}), \qquad 1143$$

taking $\frac{0}{0} \coloneqq 0$. 1144

B.2.2 Pairwise Ranking

We prompt the LLM judge to express its preference on a *K*-point Likert scale: [>, <] (Likert-2), [>, =, <] (Likert-3), or $[\gg, >, =, <, \ll]$ (Likert-5). Assigning the symbols $[\gg, >, =, <, \ll]$ the numerical values [2, 1, 0, -1, -2], the methods for pairwise ranking then follow those above for pairwise scoring. We remark that the 'mode' and 'mean' for pairwise scoring and pairwise ranking in Table 1 are with post-aggregation and preaggregation, respectively.

B.3 Listwise Methods

The listwise methods are introduced in Section 6.1.

C Prompts

We present representative example prompts to illustrate the different settings. The prompts are adapted from MT-Bench (Zheng et al., 2023b). Auxiliary modifications are not shown, such as the prompt for second-turn evaluation in MT-Bench.

C.1 Judgment Extraction Details

To identify the token position containing the judgment, we use the specified format when available

1129

1137

1138

1139

1140

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

(e.g. "Rating A: {rating_a}." in pairwise scor-1167 ing). Otherwise, we use the latest token position 1168 with more than 0.5 total probability assigned to 1169 judgment tokens. If no valid token is found, we de-1170 fault the judgment to the minimum score of 1 in the 1171 scoring setting, and to a tie in the ranking setting. 1172 (For Nectar experiments, we exclude instances with 1173 invalid silver-label judgments.) 1174

> For the local models (Llama-3.1-8B, Mistral-7B, Prometheus-2-7B) in no-CoT prompting, we force a prefix of the assistant's response (e.g. "Rating A: ") and use a single output token as the judgment token position.

1180

1175

1176

1177

1178

1179

C.2 Pointwise Prompts

System prompt for pointwise scoring (CoT, K = 9)

Please act as an impartial judge and

- \rightarrow evaluate the quality of the response
- $\, \hookrightarrow \ \ \, \text{provided by an AI assistant to the user}$
- \rightarrow prompt displayed below. Your
- $\ \ \, \rightarrow \quad evaluation \ should \ consider \ factors$
- \hookrightarrow such as the helpfulness, relevance,
- $\ \ \, \rightarrow \quad accuracy, \, depth, \, creativity, \, level \, of$
- \hookrightarrow detail, and ethicality of the response.
- \rightarrow Begin your evaluation by providing a
- \leftrightarrow short explanation. Be as objective as
- \rightarrow possible. After providing your
- \rightarrow explanation, please rate the response
- $~~ \rightarrow ~~ with ~an~integer~score~from~1~to~9,$
- \leftrightarrow without further explanation.

1181

System prompt for pointwise scoring (no-CoT, K = 9)

Please act as an impartial judge and

- \rightarrow evaluate the quality of the response
- \rightarrow provided by an AI assistant to the user
- \rightarrow prompt displayed below. Your
- \leftrightarrow evaluation should consider factors
- \rightarrow such as the helpfulness, relevance,
- \rightarrow accuracy, depth, creativity, level of
- \rightarrow detail, and ethicality of the response.
- \rightarrow Be as objective as possible. Please rate
- $\, \hookrightarrow \,$ the response with an integer score
- \hookrightarrow from 1 to 9, without further
- \rightarrow explanation.

User prompt for pointwise judgment

[User Prompt] {User Prompt} [End User Prompt]

[Start of Assistant's Answer]
{Assistant's Answer}
[End of Assistant's Answer]

1183

C.3 Pairwise Prompts

1184

System prompt for pairwise scoring (CoT, K = 9)

Please act as an impartial judge and

- \rightarrow evaluate the quality of the responses
- $\, \hookrightarrow \quad \text{provided by two AI assistants to the} \,$
- $\hookrightarrow \quad$ user prompt displayed below. Your
- $\, \hookrightarrow \quad evaluation \ should \ consider \ factors$
- $\, \hookrightarrow \quad \text{such as the helpfulness, relevance,} \,$
- $\, \hookrightarrow \,$ accuracy, depth, creativity, level of
- \rightarrow detail, and ethicality of their responses.
- \rightarrow Begin your evaluation by comparing
- \rightarrow the two responses and provide a short \rightarrow explanation. Avoid any position biases
- \rightarrow explanation. Avoid any position biases \rightarrow and ensure that the order in which the
- \rightarrow and ensure that the order in which th \rightarrow responses were presented does not
- \rightarrow influence your decision. Do not allow
- \rightarrow the length of the responses to
- \rightarrow influence your evaluation. Do not
- \rightarrow favor certain names of the assistants.
- \rightarrow Be as objective as possible. After
- \rightarrow providing your explanation, output
- \rightarrow your final verdict by strictly following
- \rightarrow this format: "Rating A: {rating_a}.
- \rightarrow Rating B: {rating_b}.", where
- \rightarrow "{rating_a}" and "{rating_b}" are

For pairwise ranking with the local models, we

use a different prompt from the one below. We

found that they would often fail to include the

braces specified in the judgment format, so we

omit them when prompting these models.

 \rightarrow integer scores from 1 to 9.

1186 1187

1185

1188 1189

1190

Please act as an impartial judge and

- $\hookrightarrow \quad evaluate \ the \ quality \ of \ the \ responses$
- \hookrightarrow provided by two AI assistants to the
- $\hookrightarrow \quad \text{user prompt displayed below. You}$
- \hookrightarrow should choose the assistant that
- $\hookrightarrow \quad follows \ the \ user's \ instructions \ and$
- $\, \hookrightarrow \,$ answers the user's question better.
- \hookrightarrow Your evaluation should consider
- \hookrightarrow factors such as the helpfulness,
- $\, \hookrightarrow \,$ relevance, accuracy, depth, creativity,
- $\hookrightarrow \quad \text{level of detail, and ethicality of their}$
- $\, \hookrightarrow \,$ $\,$ responses. Begin your evaluation by
- $\, \hookrightarrow \ \ \, comparing \ the \ two \ responses \ and$
- $\hookrightarrow \quad \mbox{provide a short explanation. Avoid any}$
- \hookrightarrow position biases and ensure that the
- \hookrightarrow order in which the responses were
- \hookrightarrow presented does not influence your
- \rightarrow decision. Do not allow the length of
- $\hookrightarrow \quad \text{the responses to influence your}$
- $\, \hookrightarrow \,$ evaluation. Do not favor certain names
- $\, \hookrightarrow \,$ of the assistants. Be as objective as
- $\ \ \, \rightarrow \quad possible. \ After \ providing \ your$
- \leftrightarrow explanation, output your final verdict
- \rightarrow by strictly following this format:
- \rightarrow "[[>>]]" if assistant A is significantly
- \rightarrow better, "[[>]]" if assistant A is slightly
- \rightarrow better, "[[=]]" for a tie, "[[<]]" if
- \rightarrow assistant B is slightly better, and
- \rightarrow "[[<<]]" if assistant B is significantly

 \rightarrow better.

1191

User prompt for pairwise judgment

[User Prompt] {User Prompt} [End User Prompt]

[Start of Assistant A's Answer] {Assistant A's Answer} [End of Assistant A's Answer]

[Start of Assistant B's Answer]
{Assistant B's Answer}
[End of Assistant B's Answer]

C.4 Listwise Prompts

The listwise prompts are adapted from Nectar (Zhu1194et al., 2024).1195

System prompt for listwise judgment (N = 7), with intermediate pairwise preferences

We are interested in ranking different large

- \leftrightarrow language model chat completions to a
- \hookrightarrow conversation. Please act as an
- \rightarrow impartial judge and evaluate the
- \rightarrow quality of the completions provided by
- \leftrightarrow the 7 AI assistants. Your evaluation
- \rightarrow should consider factors such as the
- \rightarrow helpfulness, relevance, accuracy,
- \rightarrow depth, creativity, level of detail, and
- \rightarrow ethicality of their responses.

After the conversation and assistant

- \rightarrow responses, the section "PAIRWISE
- → EVALUATION ORDER" will specify
- \rightarrow the order in which to perform pairwise
- \rightarrow comparisons. Output an array in
- \rightarrow which, for each pairwise comparison,
- \rightarrow vou choose the letter of the better
- \rightarrow response, or '=' for a tie. The array
- \rightarrow should be comma-separated and
- \leftrightarrow enclosed in double square brackets.

Then, considering these pairwise rankings,

- \rightarrow please rank all 7 responses from best
- \rightarrow to worst (breaking ties randomly),
- \rightarrow strictly in the following format: [[_, _,
- \rightarrow _, _, _, _, _]] where '_' contains an
- \rightarrow assistant's letter name.

Avoid any position biases and ensure that

- \rightarrow the order in which the responses were
- \rightarrow presented does not influence your
- \rightarrow decision. Do not allow the length of
- \hookrightarrow the responses to influence your
- $\, \hookrightarrow \,$ evaluation. Do not favor certain names
- \rightarrow of the assistants. Be as objective as
- \rightarrow possible.

System prompt for listwise judgment (N = 7), without intermediate pairwise preferences

We are interested in ranking different large

- \rightarrow language model chat completions to a
- \hookrightarrow conversation. Please act as an
- \rightarrow impartial judge and evaluate the
- \rightarrow quality of the completions provided by
- \rightarrow the 7 AI assistants. Your evaluation
- \rightarrow should consider factors such as the
- \rightarrow helpfulness, relevance, accuracy,
- \rightarrow depth, creativity, level of detail, and
- \rightarrow ethicality of their responses.

Please rank all 7 responses from best to

- \rightarrow worst (breaking ties randomly),
- \hookrightarrow strictly in the following format: [[_, _,
- \rightarrow _, _, _, _, _]] where '_' contains an
- \hookrightarrow assistant's letter name.

Avoid any position biases and ensure that

- \rightarrow the order in which the responses were
- \rightarrow presented does not influence your
- \rightarrow decision. Do not allow the length of
- \hookrightarrow the responses to influence your
- $\, \hookrightarrow \quad evaluation. \ Do \ not \ favor \ certain \ names$
- \hookrightarrow of the assistants. Be as objective as
- $\, \hookrightarrow \ \ possible.$

1197

User prompt for listwise judgment (N = 7). The presentation order is randomized. The pairwise evaluation order is randomized every instance for the prompt with intermediate pairwise preferences, and omitted for the prompt without intermediate pairwise preferences.

[CONVERSATION START] {Conversation} [CONVERSATION END]

[MODEL A RESPONSE START] {Model A's response} [MODEL A RESPONSE END]

[MODEL B RESPONSE START] {Model B's response} [MODEL B RESPONSE END]

[MODEL C RESPONSE START] {Model C's response} [MODEL C RESPONSE END]

[MODEL D RESPONSE START] {Model D's response} [MODEL D RESPONSE END]

[MODEL E RESPONSE START] {Model E's response} [MODEL E RESPONSE END]

[MODEL F RESPONSE START] {Model F's response} [MODEL F RESPONSE END]

[MODEL G RESPONSE START] {Model G's response} [MODEL G RESPONSE END]

PAIRWISE EVALUATION ORDER: [(G, \rightarrow C), (B, G), (C, D), (A, E), (G, A), (A, \rightarrow D), (B, A), (B, E), (B, F), (A, C), (E, \rightarrow C), (E, F), (B, D), (F, A), (G, E), (F, \rightarrow C), (F, D), (C, B), (F, G), (D, G), (E, \rightarrow D)]

D Datasets

17

RewardBench (Lambert et al., 2024) is a reward model benchmark spanning chat, reasoning, and 1198

1199

safety. Each instance consists of a prompt, a cho-1202 sen response, and a rejected response, all manually 1203 verified. The dataset categories are Chat, with 358 1204 instances sourced from AlpacaEval (Li et al., 2023) 1205 and MT-Bench (Zheng et al., 2023b); Chat Hard, with 456 instances sourced from MT-Bench and 1207 LLMBar (Zeng et al., 2023); Safety, with 740 in-1208 stances sourced from XSTest (Röttger et al., 2023), 1209 Do-Not-Answer (Wang et al., 2023), and original 1210 data; and Reasoning, with 1431 instances sourced 1211 from PRM800k (Lightman et al., 2023) and Hu-1212 manEvalPack (Muennighoff et al., 2023). Except 1213 for excluding the prior sets category, we follow the 1214 original work and compute the final score as the 1215 average of the category scores. 1216

1217

1218

1219

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

MT-Bench (Zheng et al., 2023b) is a dataset of multi-turn questions spanning writing, roleplay, extraction, reasoning, math, coding, knowledge I (STEM), and knowledge II (humanities/social science). There are 3,355 (prompt, model pair, human judge, turn) tuples, 1,814 unique (prompt, model pair, turn) tuples, and 80 unique prompts each with two turns of interaction. To evaluate accuracy, we use the 1,132 instances with unanimous non-tie human judgments. To evaluate MSE, we use all 1,814 instances and set the label of an instance to the average of the human judgments, where a 0 or 1 represents the evaluated winner, and a 0.5 represents a tie.

Nectar (Zhu et al., 2024) is a dataset of 183k prompts each with 7 model responses. The prompts are sourced from Anthropic-HH (Bai et al., 2022), LMSYS-Chat-1M (Zheng et al., 2023a), UltraFeedback (Cui et al., 2023), and ShareGPT. We use a random subset of size 1,000.

RM-Bench (Liu et al., 2024c) is a reward model benchmark focusing on sensitivity to subtle content differences and resistance to style biases. There are 1,327 instances spanning chat, code, math, and safety. Similar to RewardBench, we follow the original work and average the 4 category scores. For each prompt, there are 3 pairs of (chosen, rejected) responses, where each pair is written with a particular style regarding concision and whether formatted as plain text or markdown.

The HelpSteer2 dataset (Wang et al., 2024b) contains multiple human ratings on a 0-4 scale for five attributes (helpfulness, correctness, coherence, complexity, verbosity) for each (prompt, response) instance. We use a random subset of size 1,000.

D.1 Listwise Evaluation

For the listwise setting, we use the same evaluation setup as with the pointwise and pairwise setting.⁴ We concern ourselves with agreement at the pair level rather than the list level because pairwise preferences are sufficient to produce a total order, such as by choosing the maximum likelihood order (Liu et al., 2024e; Liusie et al., 2024) or with graphtheoretic methods (Tideman, 1987; Schulze, 2011; Li et al., 2024c). Thus, pairwise preferences are an adequate unit at which to measure agreement, and the aggregation into a total order may be modularized away for experimental simplicity. 1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1289

1290

1291

1292

1294

1295

1296

1297

1298

To compute accuracy on Nectar with silver labels (Section 6.2), we take the sign of the silver label as the silver label for accuracy.

E Additional Results

Table 10 is an expanded version of Table 1, providing a subset breakdown of RewardBench. We observe particularly large gains for pointwise scoring on the Reasoning subset, e.g. absolute +7.7% and +17.1% for GPT-40 and Llama-3.1-8B.

Tables 11 (K = 9) and 14 (K = 99) show pointwise results over methods (expanded versions of Tables 4 and 6). Tables 12 and 13 show expanded tie analyses on RewardBench (simplified in Table 5) and MT-Bench.

Table 15 shows pairwise ranking results over methods, extending Table 7. Table 16 compares the Likert scales used for pairwise ranking, extending Table 8. In Table 16, the most calibrated setting on MT-Bench is (GPT-40) Likert-5 no-CoT, achieving a 31% lower MSE than the most accurate setting, Likert-2 CoT, suggesting that a finer granularity has potential to improve calibration (Liu et al., 2024b). With GPT-40 in Tables 15 and 16, MSE is always minimized with no-CoT, highlighting the discord between CoT's sharpening effect and calibration. This result is in line with AlpacaEval (Dubois et al., 2024), which uses no-CoT and judgment probabilities, but deviating from WB-Reward and Arena-Hard-Auto (Lin et al., 2024; Li et al., 2024b), which use CoT and decoded judgments.

E.1 DeepSeek-V3 Results

We provide partial results for DeepSeek-V3 (DeepSeek-AI et al., 2025), a model of comparable size to GPT-40. Tables 17 and 18 contain pointwise

⁴This means that our MT-Bench results are directly comparable across settings.

Model	Setting	Method		RewardBench				
	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		Chat	Chat Hard	Safety	Reasoning	Total	
	• ,	mode	95.8, 89.7	76.0, 77.4	89.3, 88.5	79.5, 80.3	85.1, 84.0	81.9, 80.5
	point score	mean	<b>97.1</b> , 94.3	75.2, <b>79.8</b>	<b>90.3</b> , <u>89.7</u>	<u>87.0</u> , <b>88.0</b>	<u>87.4</u> , <b>88.0</b>	<b>83.6</b> , <u>83.2</u>
CDT 40	noir sooro	mode	<u>97.3</u> , <b>97.9</b>	69.0, <u>70.7</u>	<u>89.1</u> , <b>89.5</b>	91.3, 91.3	86.7, 87.4	86.2, 86.5
OF 1-40	pair score	mean	<u>97.2</u> , <u>97.8</u>	<u>69.7</u> , <b>70.8</b>	<b>89.5</b> , <b>89.5</b>	<u>91.9</u> , <b>92.4</b>	<u>87.1</u> , <b>87.6</b>	<u>86.3</u> , <b>86.8</b>
	nair rank	mode	96.9, <u>97.6</u>	76.4, <u>79.1</u>	89.0, <b>90.9</b>	91.4, 91.3	88.4, 89.7	86.3, 85.6
	pair raik	mean	96.2, <b>98.3</b>	76.6, <b>79.4</b>	88.5, <u>90.8</u>	<u>93.0</u> , <b>93.6</b>	88.6, <b>90.5</b>	<b>87.3</b> , 85.9
	point score	mode	83.8, 87.6	<u>57.6, 58.0</u>	76.2, 78.2	60.8, 64.8	69.6, 72.2	74.9, 71.9
	point score	mean	89.0, <b>95.8</b>	<u>58.6</u> , <b>58.8</b>	73.0, <b>80.8</b>	70.2, <b>81.9</b>	72.7, <b>79.3</b>	78.7, <b>81.5</b>
Llama_3_1_8B	pair score	mode	92.0, 94.3	45.4, 45.4	69.5, <u>78.8</u>	79.9, 82.4	71.7, 75.2	<b>82.6</b> , <u>82.4</u>
Liama-5.1-6D		mean	92.6, <b>95.8</b>	<u>44.6, 45.0</u>	69.3, <b>78.9</b>	81.7, <b>87.6</b>	72.1, <b>76.8</b>	<u>82.3</u> , 81.2
	pair rank	mode	76.7, 65.2	<b>52.3</b> , 48.1	71.0, 66.4	75.6, 55.8	68.9, 58.9	76.2, 63.0
		mean	<u>90.5</u> , <b>93.0</b>	<u>50.0</u> , 44.1	<b>78.1</b> , 72.7	<b>78.3</b> , 64.6	<b>74.2</b> , 68.6	<b>80.0</b> , 76.5
	maint sagera	mode	52.4, 66.2	<u>51.5, 50.5</u>	<b>79.9</b> , 75.7	57.8, 58.4	60.4, 62.7	59.5, 66.2
	point score	mean	54.5, <b>82.1</b>	<b>53.5</b> , <u>49.1</u>	<b>79.9</b> , <u>79.6</u>	67.2, <b>77.5</b>	63.8, <b>72.1</b>	62.6, <b>74.0</b>
Mistral 7B	pair score	mode	87.6, <u>89.9</u>	<u>40.2, 40.4</u>	<u>74.0, 73.0</u>	67.4, 72.4	67.3, 68.9	<u>79.3, 79.8</u>
Wilsual-7D	pair score	mean	<u>89.2</u> , <b>91.1</b>	<b>41.2</b> , <u>39.3</u>	<b>74.1</b> , <u>73.4</u>	67.8, <b>80.2</b>	68.1, <b>71.0</b>	<u>80.0</u> , <b>80.4</b>
	noir ronk	mode	51.0, 51.5	<b>51.0</b> , 46.2	62.2, 66.8	61.0, 50.8	56.3, 53.8	51.5, 51.5
	pair raik	mean	<u>79.5</u> , <b>81.7</b>	39.3, 36.3	<b>73.1</b> , 67.7	<b>63.8</b> , 50.6	<b>63.9</b> , 59.1	<b>73.5</b> , 65.5
	point score	mode	81.3, 81.7	50.5, 50.8	65.9, 73.4	59.2, 58.2	64.3, 66.0	72.5, 73.5
	point score	mean	82.4, <b>92.2</b>	48.9, <b>54.4</b>	65.7, <b>76.6</b>	61.3, <b>77.6</b>	64.6, <b>75.2</b>	72.1, <b>81.6</b>
Prometheus 2 7B	pair score	mode	<u>91.2, 92.0</u>	<b>44.1</b> , <u>43.6</u>	<b>75.9</b> , 69.4	72.7, 69.6	<b>71.0</b> , 68.7	78.4, <u>80.8</u>
Prometneus-2-/B	pair score	mean	<u>91.3</u> , <b>93.0</b>	42.7, <u>43.0</u>	74.9, 72.0	<u>73.0</u> , <b>75.1</b>	70.5, <u>70.8</u>	78.3, <b>80.9</b>
	nair rank	mode	55.6, 45.4	<b>5</b> 1.0, <u>50.0</u>	66.6, 49.7	65.3, 47.8	59.6, 48.2	51.5, 43.0
	pair rank	mean	<b>90.5</b> , 45.0	44.3, <u>50.7</u>	<b>74.2</b> , 55.5	<b>69.8</b> , 44.1	<b>69.7</b> , 48.8	<b>75.4</b> , 33.4

Table 10: Mode vs. mean and CoT vs. no-CoT (comma-separated) accuracy results (%). Expanded version of Table 1.

and listwise results, respectively. The trends for DeepSeek-V3 match those of GPT-40.

#### F Analysis

1299

1300

1301

1303

1304

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318 1319

1320

1321

1322

1323

#### F.1 Heterogenous Preferences

We investigate whether LLM judges can represent pluralistically aligned preferences (i.e. reflect diverse human opinions) (Sorensen et al., 2024; Siththaranjan et al., 2023; Kumar et al., 2024) through their judgment distribution, without explicit training or prompting.

#### F.1.1 Multimodality

We begin by quantifying the degree of multimodality in the judgment distributions. An implicit assumption behind the conventional method of using the mode judgment is that the judgment distribution is unimodal and thus the mode is a representative judgment. However, in cases where humans disagree, we would like LLM judges to reflect the heterogeneity in the human population with a multimodal distribution.

We quantify multimodality as the minimum amount of probability mass that must be added to make an unnormalized unimodal distribution, divided by the total mass of the unnormalized unimodal distribution to obtain a value in [0, 1), where a distribution is unimodal if the probability mass function is non-decreasing and then nonincreasing. For example, if the judgment distribution is [0.5, 0.2, 0.3], the minimum additional mass is 0.1 to obtain the unimodal distribution [0.5, 0.3, 0.3] with total mass 1.1, so we compute the multimodality as  $0.1/1.1 \approx 0.091$ . 1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

Table 19 presents the results. We find that more granularity leads to more multimodality (note that K = 2 always has multimodality 0), and no-CoT is more multimodal than CoT. The case of extreme multimodality for pointwise scoring K = 99 can be largely attributed to token bias (Lovering et al., 2024; Shaikh et al., 2024). For example, GPT-40 K = 99 CoT on MT-Bench assigns on average 0.036 probability to a single token that is a multiple of 5, but only 0.002 to a single token that differs by 1 from one of those multiples of 5.

#### F.1.2 Annotator Disagreement

We next examine whether human annotator dis-<br/>agreement is correlated with the uncertainty in the<br/>LLM's judgment distribution. On datasets with<br/>multiple human judgments per instance, we com-<br/>pute Spearman's  $\rho$  between the standard deviation<br/>of the human judgments and that of the LLM's<br/>judgment distribution.1343<br/>13441345<br/>1346<br/>13471346<br/>1347

Model	Method	Reward	dBench	MT-E	Bench
		Acc ↑	MSE↓	Acc ↑	$MSE\downarrow$
GPT-4o	MODE MEAN [MEAN] MEDI 1P RAM QT PS	85.1, 84.0 87.4, 88.0 85.1, 85.2 85.0, 84.6 84.8, 84.3 87.4, 88.4 87.4, 87.9 87.4, 87.8	.116, .118 .099, .102 .116, .109 .116, .112 .120, .116 .099, .100 .107, .096 .106, .096	81.9, 80.5 83.6, <u>83.2</u> 82.0, 80.2 82.0, 80.2 82.6, 81.0 <b>83.9</b> , <u>83.4</u> 83.5, <u>83.2</u> 83.5, <u>83.3</u>	.152, .145 .115, <u>.097</u> .150, .146 .150, .142 .141, .138 .115, <b>.096</b> .139, .118 .136, .103
Llama -3.1-8B	MODE MEAN [MEAN] MEDI 1P RAM QT PS	69.6, 72.2 72.7, 79.3 70.1, 75.0 69.8, 73.6 70.2, 76.0 72.7, <b>79.9</b> 72.8, 79.0 72.8, 78.9	.237, .192 .198, .155 .238, .186 .238, .191 .238, .183 .200, <b>.152</b> .220, .164 .216, .161	74.9, 71.9 78.7, <b>81.5</b> 75.7, 75.0 75.2, 73.9 76.8, 79.2 78.8, <u>81.4</u> 78.7, 81.1 78.6, <u>81.4</u>	.177, .142 .129, .104 .172, .145 .176, .142 .172, .147 .130, <b>.102</b> .154, .116 .149, .110

Table 11: Pointwise results over methods. Comma-separated values are with and without CoT (expanded version of Table 4). Text styling follows Table 1.

Model	Method	Tie rate		MEAN's accuracy		Non-tie accuracy $\Delta \uparrow$	
		K = 9	K = 99	K = 9	K = 99	K = 9	K = 99
	MODE	.13, .17	.09, .20	64, 72	61, 73	+0.0, -0.1	+0.0, -0.3
CDT 4a	[MEAN]	.13, .16	.02, .03	65, 67	61, 53	+0.0, +0.0	+0.0, -0.0
OF 1-40	MEDI	.13, .17	.06, .09	65, 70	58, 62	-0.0, +0.0	-0.0, -0.1
	1 P	.13, .16	.05, .08	66, 66	58,60	+0.1, +0.1	+0.0, +0.5
	MODE	.27, .35	.18, .24	60, 69	63, 70	+0.2, -0.2	+0.2, -2.0
Llama	[MEAN]	.25, .26	.07, .07	58, 64	56, 61	+0.0, +0.0	+0.0, +0.0
-3.1-8B	MEDI	.26, .29	.11, .11	60, 67	59, 67	+0.1, -0.5	-0.2, -0.6
	1 P	.24, .23	.08, .08	61, 65	55, 57	+0.3, +0.6	+0.8, +0.1

Table 12: Tie analysis for discrete pointwise methods on RewardBench (expanded version of Table 5). Tie rate is the proportion of instances where the method predicts a tie, over which we report MEAN's accuracy (%); excess of 50% or 75% indicates room for improving accuracy or MSE, respectively. Non-tie accuracy  $\Delta$  (%) is the method's accuracy minus MEAN's accuracy over the non-tie instances. Comma-separated values are with and without CoT. We find that the mode has the most ties, the highest MEAN accuracy, and the lowest non-tie accuracy delta (i.e. poor recall without better precision), especially for no-CoT K = 99.

For MT-Bench, we take the 961 instances with multiple human judgments. Table 20 reports weak correlation in all settings except no correlation in pairwise ranking with Llama-3.1-8B. Remarkably, *pointwise* score distributions encode sufficient information to predict if humans will disagree on a *pairwise* comparison of the texts.

1350

1351

1352

1353

1354

1355

1356

1358

1359

1360

1361

1362

1363

1364

1366

The HelpSteer2 dataset (Wang et al., 2024b) contains multiple human ratings on a 0-4 scale for five attributes for each (prompt, response) instance. We use a random subset of size 1,000. We prompt with the provided annotation guidelines and have the model rate all attributes in a single run. Table 21 reports weak correlation on helpfulness, correctness, and coherence but no correlation on complexity and verbosity. We suspected this to be due to that conditioning on the earlier attributes' scores may reduce uncertainty for the later attributes (Stureborg et al., 2024; Hashemi et al., 2024), but we found that the average standard deviation is similar across attributes for both LLM and human judgments. 1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

#### F.1.3 Pluralistic Alignment

We finally evaluate the alignment between predicted judgment distributions and human judgment distributions. We quantify the distance between two distributions  $\mu$  and  $\nu$  with the Wasserstein *p*distance for  $p \in \{1, 2\}$ :

$$W_p(\mu,\nu) = \inf_{\gamma \in \Gamma(\mu,\nu)} \left( \mathop{\mathbb{E}}_{(x,y) \sim \gamma} |x-y|^p \right)^{\frac{1}{p}}, \quad (1)$$
 1377

where  $\Gamma(\mu, \nu)$  is the set of couplings of  $\mu$  and  $\nu$ . A higher p more heavily punishes large point distances |x - y|. We scale the judgment spaces to

Model	Method	Tie rate		MEAN's accuracy		Non-tie accuracy $\Delta \uparrow$	
		K = 9	K = 99	K = 9	K = 99	K = 9	K = 99
	MODE	.13, .21	.08, .22	62, 62	64, 67	+0.0, -0.2	+0.1, -0.7
GPT-40	[MEAN]	.13, .21	.02, .04	61,64	42, 48	+0.0, +0.0	+0.0, +0.0
	MEDI	.13, .22	.05, .11	61,64	64, 55	+0.0, +0.1	+0.0, -0.6
	1 P	.14, .19	.06, .09	56, 62	67, 56	-0.1, +0.1	+0.3, +0.7
	MODE	.25, .45	.14, .26	65, 71	61, 65	-0.1, -1.0	-0.4, -3.0
Llama	[MEAN]	.24, .36	.06, .09	63, 68	49, 55	+0.0, +0.0	+0.0, -0.1
-3.1-8B	MEDI	.25, .40	.10, .18	65, 69	54, 58	+0.0, -0.3	-0.4, +0.5
	1P	.20, .23	.07, .07	60, 59	53, 50	+0.1, -0.3	+1.8, +0.3

Table 13: Tie analysis for discrete pointwise methods on MT-Bench, mirroring Table 12.

Model	Method	Rewa	rdBench	MT-	Bench
		Acc ↑	MSE ↓	Acc ↑	MSE ↓
	MODE	86.1+1.0, 81.7-2.3	$.118_{+.002}, .134_{+.016}$	83.8+1.9, 78.4_2.1	$.152_{+.000}, .158_{+.013}$
	MEAN	87.4 _{+0.0} , 86.7 _{-1.3}	<b>.097</b> ₀₀₂ , .108 _{+.006}	$84.8_{+1.2}, 82.9_{-0.3}$	$.105_{010}, .099_{+.002}$
	[MEAN]	$87.0_{+1.9}, 86.5_{+1.3}$	$.124_{+.008}, .127_{+.018}$	<b>85.1</b> _{+3.1} , 82.7 _{+2.5}	$.167_{+.017}, .182_{+.036}$
CDT 4a	MEDI	86.7+1.7, 85.2+0.6	$.119_{+.003}, .126_{+.014}$	84.1+2.1, 81.5+1.3	$.160_{+.010}, .170_{+.028}$
GP1-40	1 P	86.6+1.8, 86.4+2.1	$.121_{+.001}, .116_{+.000}$	$84.2_{+1.6}, 82.7_{+1.7}$	$.159_{+.018}, .165_{+.027}$
	RAM	87.1 _{-0.3} , <u>86.7</u> _{-1.7}	$.098_{001}, .104_{+.004}$	<b>85.1</b> _{+1.2} , 83.0 _{-0.4}	.106009, .098+.002
	QT	87.3_0.1, 86.6_1.3	$.112_{+.005}, .114_{+.018}$	$84.8_{+1.3}, 82.7_{-0.5}$	$.149_{+.010}, .147_{+.029}$
	PS	$\overline{87.3}_{-0.1}, \overline{86.6}_{-1.2}$	$.105_{001}, .105_{+.009}$	$\overline{84.8}_{+1.3}, 82.4_{-0.9}$	$.130_{006}, .107_{+.004}$
	MODE	73.4+3.8, 72.0-0.2	.222015, .221+.029	77.3+2.4, 75.1+3.2	$.191_{+.014}, .169_{+.027}$
	MEAN	$75.9_{+3.2}, \underline{79.3}_{+0.0}$	$.183_{015}, .156_{+.001}$	$79.3_{+0.6}, \underline{81.3}_{-0.2}$	.125004, .103001
	[MEAN]	75.3+5.2, 78.5+3.5	.229009, .198+.012	79.3+3.6, 80.7+5.7	$.201_{+.029}, .180_{+.035}$
Llama	MEDI	$74.4_{+4.6}, 76.5_{+2.9}$	.228010, .207+.016	$78.4_{+3.2}, 80.1_{+6.2}$	$.198_{+.022}, .161_{+.019}$
-3.1-8B	1 P	$76.2_{+6.0}, 78.5_{+2.5}$	$.218_{020}, .195_{+.012}$	$80.6_{+3.8}, 81.5_{+2.3}$	$.187_{+.015}, .177_{+.030}$
	RAM	76.1 _{+3.4} , <b>79.7</b> _{-0.2}	.179021, .152+.000	$79.7_{+0.9}, 81.1_{-0.3}$	.123007, .102+.000
	QT	$75.7_{+2.9}, 78.7_{-0.3}$	$.214_{006}, .177_{+.013}$	$78.8_{\pm 0.1}, \overline{81.3}_{\pm 0.2}$	$.179_{+.025}, .143_{+.027}$
	PS	75.7+2.9, 78.6-0.3	.203013, .163+.002	78.6+0.0, <b>81.8</b> +0.4	$.151_{+.002}, .111_{+.001}$

Table 14: Pointwise results over methods (K = 99). Comma-separated values are with and without CoT (expanded version of Table 6). Subscripts denote change from K = 9 (Table 11). Text styling follows Table 1.

[0,1] so that  $W_p(\mu,\nu) \in [0,1]$ .

1381

1382

1383

1384

1386

1387

1388

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

As baselines, we consider deterministic distributions that place probability 1 on a measure of central tendency.

Table 22 shows that using a distributional prediction has little success in improving alignment with the MT-Bench human pairwise preferences, but Table 23 shows success for HelpSteer2 human pointwise scores.

We also experimented with the HelpSteer2-Preference dataset, prompting with the provided annotation guidelines (Wang et al., 2024a). However, we found severe position bias in our experiments with GPT-40 and Llama-3.1-8B (no-CoT). The analysis showed no correlation between predicted distribution variance and annotator disagreement, and poor pluralistic alignment compared to the deterministic baselines.

## F.2 Sensitivity to Score Granularity

Adopting the view that LLMs latently encode a continuous distribution but output a discretization of it (Gillman et al., 2024), we analyze how faithfully 1402 functions of the (latent) continuous distribution can 1403 be approximated by those functions computed on 1404 the (observed) discretization. For practical inter-1405 est, this manifests as robustness to the choice of 1406 K, with convergence in distribution to the continu-1407 ous distribution as  $K \to \infty$ . Thus, independently 1408 of the "principledness" of certain functions of a 1409 ground-truth continuous distribution, it is appropri-1410 ate to examine the effect of discretization on our 1411 ability to approximate them to begin with. Our 1412 theoretical result is stated in Proposition 1 (see 1413 Appendix G.1 for full statement, proof, and discus-1414 sion). 1415

**Proposition 1.** Among the discrete methods in Ta-<br/>ble 3, MODE computed on continuous distributions1416may fail to be approximated by the same function<br/>computed on their discretizations, even under regu-<br/>larity conditions. Meanwhile, [MEAN], MEDI, and<br/>142014201P admit an approximation error bound.1421

We empirically assess the robustness to K of

Model	Center	Agg.	Reward	lBench	MT-Bench		
		Time	Acc ↑	$MSE\downarrow$	$Acc \uparrow$	$MSE\downarrow$	
	mada	post	88.1, 89.3	.099, <b>.090</b>	<u>86.1</u> , 84.9	<b>.139</b> , <u>.142</u>	
GPT-40	mode	pre	88.4, <b>90.3</b>	.112, .094	<b>86.5</b> , 85.2	.154, .154	
	median	post	88.1, 89.3	.099, <b>.091</b>	<u>86.1</u> , 84.9	<b>.138</b> , <u>.142</u>	
	methan	pre	88.4, <b>90.0</b>	.111, <u>.094</u>	<b>86.6</b> , 85.4	.153, .146	
	mean	post	88.9, <b>90.4</b>	.098, <b>.077</b>	<b>86.5</b> , 85.4	.132, .100	
		pre	88.9, <b>90.4</b>	.098, <u>.078</u>	<b>86.6</b> , 85.4	.132, <b>.097</b>	
	mode	post	56.7, 52.4	<b>.240</b> , .279	57.5, 53.4	.192, <b>.176</b>	
	moue	pre	<b>73.1</b> , 66.1	.265, .337	<b>78.1</b> , 70.9	.222, .268	
Llama	median	post	56.8, 52.5	<b>.240</b> , .279	57.5, 53.5	.192, <b>.176</b>	
-3.1-8B	methan	pre	<b>72.9</b> , 65.3	.261, .319	<b>78.0</b> , 69.1	.218, .238	
	mean	post	<u>73.2</u> , 65.6	<b>.207</b> , .229	<b>78.2</b> , 70.5	<b>.144</b> , <u>.146</u>	
	mean	pre	<b>73.2</b> , 66.3	.222, .240	<u>78.1</u> , 70.8	.155, .155	

Table 15: Pairwise ranking results over methods, using Likert-3 (expanded version of Table 7). Comma-separated values are with and without CoT. Text styling follows Table 1.

Model	K	Reward	dBench	MT-Bench		
		Acc $\uparrow$	$MSE\downarrow$	Acc $\uparrow$	$MSE\downarrow$	
	2	88.6, <b>90.5</b>	.094, <b>.077</b>	<b>87.3</b> , 85.9	.136, .101	
GPT-40	3	88.9, <u>90.4</u>	.098, <b>.078</b>	<u>86.6</u> , 85.4	.132, .097	
	5	88.8, 89.5	.099, .106	84.7, 85.8	.129, <b>.087</b>	
Llama	2	<b>74.2</b> , 68.6	<b>.187</b> , .214	<b>80.0</b> , 76.5	<b>.126</b> , .135	
-3.1-8B	3	<u>73.2</u> , 66.3	.222, .240	78.1, 70.8	.155, .155	
	5	70.0, 58.5	.215, .234	77.1, 64.8	.142, .153	

Table 16: Pairwise ranking results over Likert-K scales, using pre-aggregation mean (expanded version of Table 8). Comma-separated values are with and without CoT. Text styling follows Table 1.

Method	Acc $\uparrow$	$MSE\downarrow$
MODE	<u>84.7</u> , 82.5	<u>0.123</u> , 0.128
MEAN	<u>84.8</u> , <u>84.2</u>	<u>0.119, 0.120</u>
[MEAN]	<u>84.7</u> , 82.7	<u>0.123</u> , 0.127
MEDI	<u>84.7</u> , 82.6	<u>0.123</u> , 0.128
1 P	84.5, 82.9	<u>0.124</u> , 0.126
RAM	<b>84.9</b> , <u>84.1</u>	<u>0.120, <b>0.118</b></u>
QT	<b>85.0</b> , 83.9	<u>0.122</u> , 0.125
PS	<b>85.0</b> , 83.9	<u>0.122</u> , 0.125

Table 17: Pointwise results with DeepSeek-V3 on RewardBench. K = 9. Comma-separated values are with and without CoT. Text styling follows Table 1.

the score distributions produced by the LLM judge as well as the functions computed on them. The former is not addressed by Proposition 1, which assumes the score distributions to be errorless discretizations and thus consistent across granularities.

#### F.2.1 Sensitivity of Score Distributions

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

For an evaluated text, let  $\mu^{K}$  denote the score distribution with granularity K, with the score space scaled to [0, 1]. We coarsify  $\mu^{99}$  into  $\hat{\mu}^{99}$ 

Space	Method	Ne	ectar	RM-Bench		
Spuee		Acc	MSE	Acc	MSE	
direct list direct list	mode mean	83.6 <b>84.0</b>	0.149 <b>0.129</b>	<b>67.8</b> <u>67.6</u>	0.322 <b>0.307</b>	

Table 18: Listwise results with DeepSeek-V3. Text styling follows Table 1.

Model	Setting	K	RewardBench	MT-Bench
GPT-4o	point score point score pair rank pair rank	9 99 3 5	.000, .008 .362, .409 .000, .018 .014, .049	.000, .012 .357, .440 .000, .019 .021, .041
Llama -3.1-8B	point score point score pair rank pair rank	9 99 3 5	.009, .040 .356, .379 .044, .091 .107, .194	.013, .025 .382, .365 .051, .081 .107, .245

Table 19: A study on multimodality (see Appendix F.1.1). Comma-separated values are with and without CoT.

by binning into 9 blocks of 11 scores. We then quantify sensitivity as the Wasserstein 1-distance  $W_1(\mu^9, \hat{\mu}^{99}) \in [0, 1]$  (Eq. 1) averaged over the pointwise instances in the dataset.

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

#### F.2.2 Sensitivity of Pointwise Methods

For a dataset  $\mathcal{D}$  of paired responses, we denote  $\mathbf{a}^K$ as the  $|\mathcal{D}|$ -length vector containing the value of a method computed on each pair using granularity K. We then quantify sensitivity as the normalized flip rate

$$FR := \frac{\|\operatorname{sgn}(\mathbf{a}^9) - \operatorname{sgn}(\mathbf{a}^{99})\|_1}{\|\operatorname{sgn}(\mathbf{a}^9)\|_1 + \|\operatorname{sgn}(\mathbf{a}^{99})\|_1} \in [0, 1].$$
(2) 1443

Model	Setting	MT-Bench
GPT-4o	point score pair score pair rank	+0.21, +0.24 +0.19, +0.27 +0.19, +0.27
Llama -3.1-8B	point score pair score pair rank	<b>+0.21</b> , <b>+0.14</b> <b>+0.20</b> , <b>+0.24</b> +0.02, -0.04

Table 20: Spearman's  $\rho$  between standard deviation of human judgments and that of LLM's judgment distribution. Comma-separated values are with and without CoT. Bold denotes significant correlation ( $\alpha = 0.01$ ). Ranking uses Likert-3; scoring uses K = 9 converted to a Likert-3 distribution  $[P(X_1 > X_2), P(X_1 = X_2), P(X_1 < X_2)].$ 

## F.2.3 Results

1444

1445

1446

1447

1448 1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1478

Table 24 presents the results on sensitivity to granularity. The discrete metrics are more sensitive than the continuous metrics. Furthermore, consistent with Proposition 1, we find that the mode is the most sensitive among the discrete methods, particularly with no-CoT.

The effect of CoT differs between the models: GPT-40 is less sensitive with CoT, and Llama-3.1-8B is less sensitive with no-CoT. Similar to Lee et al. (2024b), it would appear that although GPT-40 is a more capable judge than Llama-3.1-8B, it is not as robust to granularity (in each model's CoT/no-CoT of choice). However, this is partially because a limitation with setting K as large as 99 for GPT-40 is that no-CoT distributions tend to have high spread (Table 2), resulting in nontrivial probability mass falling outside of the top 20 tokens provided by the OpenAI API. Concretely, the average total mass on the top score tokens is 0.88/0.90 on RewardBench/MT-Bench for no-CoT, but over 0.99 for CoT.

#### F.3 Position Bias

We compare the degree of position bias (i.e. the LLM judge's sensitivity to the order in which the evaluated texts are presented (Zheng et al., 2023b)) between various settings.

1471Evaluation MetricsFor the pairwise setting1472(scoring or ranking), we measure mean absolute1473error (MAE) and mean squared error (MSE) be-1474tween the two judgments from the two orders, using1475pre-aggregation mean. Compared to MAE, MSE1476punishes a few large errors more than many small1477errors.

For the listwise setting, we measure Spearman's

 $\rho$  between the difference in the presented positions of two responses and the judgment.

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525

**Results** Tables 25 and 26 report position bias in the pairwise settings. We find that no-CoT always improves MSE, even when it hurts MAE, showing that no-CoT reduces cases of extreme position bias.

Table 27 reports listwise position bias. We find that DIRECT LIST exhibits the most position bias, consistent with Zhu et al. (2024), despite achieving the highest accuracy (Table 9). On the other hand, INTERM has the least position bias. As the intermediate pairwise preferences can be likened to CoT, this suggests that intermediate reasoning can mitigate bias in challenging judgment settings. However, since an ideal judge should be able to simultaneously maximize accuracy and minimize bias, we believe current methods have ample room for improvement.

# F.4 Transitivity

We say a comparison method  $a(\cdot, \cdot) \in [-1, 1]$  is transitive if  $a(A_1, A_2) > 0$  and  $a(A_2, A_3) \ge 0$ imply  $a(A_1, A_3) > 0$  for all triplets of texts  $(A_1, A_2, A_3)$ . For example, a score distribution comparison function that reduces to the comparison of two real numbers derived from the two score distributions independently (e.g. mode or mean) is transitive. On the other hand, QT, PS, and the pairwise ranking methods are intransitive.

Human preferences have been shown to exhibit intransitivity (Klimenko, 2015), motivating the question of whether LLM judges do so too and how this depends on the method used. Several prior works have proposed methods incorporating awareness of the intransitivity in LLM or human preferences (Liu et al., 2024e; Ethayarajh et al., 2024; Zhang et al., 2024d; Ye et al., 2024a; Hu et al., 2024; Zhang et al., 2024c; Liu et al., 2024d). We adopt the view in Liu et al. (2024e) that transitivity is generally desirable and indicative of a more capable judge, especially in the absence of a curated dataset of intransitive human preferences. Nevertheless, we remark that the ability to model intransitivity is essential to preference modeling in its full generality (Ethayarajh et al., 2024; Zhang et al., 2024d; Ye et al., 2024a), which, among pointwise methods, is achieved by QT and PS but not by mode and mean used in prior work.

⁵In every judgment space, GPT-40 tends to favor responses that are presented earlier.

Model	Helpfulness	Correctness	Coherence	Complexity	Verbosity
GPT-40	+0.24	+0.36	+0.32	+0.02	-0.01
Llama-3.1-8B	+0.14	+0.22	+0.22	-0.00	+0.01

Table 21: Spearman's  $\rho$  between standard deviation of human judgments and that of LLM's judgment distribution. HelpSteer2, no-CoT. Bold denotes significant correlation ( $\alpha = 0.01$ ).

Model	Setting	Method	$W_1$	$W_2$
		mode	.229, .246	.406, .419
	point score	mean	.229, .247	.388, .349
	1	distr	<b>.219</b> , <u>.222</u>	.395, .386
		mode	.229, .230	.419, .419
GPT-40	pair score	mean	<u>.218</u> , <b>.215</b>	.399, <b>.387</b>
		distr	<u>.220</u> , <b>.215</b>	.408, .401
		mode	.228, .226	.420, .412
	pair rank	mean	.221, .212	.396, <b>.362</b>
		distr	.215, <b>.203</b>	.405, .385
		mode	.274, .267	.438, .405
	point score	mean	.277, .267	.412, <b>.359</b>
		distr	.261, <b>.246</b>	.425, .391
Llama		mode	.268, .276	.460, .470
-3.1-8B	pair score	mean	<u>.241, .244</u>	<u>.404</u> , <b>.400</b>
-3.1-6D		distr	<b>.239</b> , <u>.243</u>	.426, .433
		mode	<b>.296</b> , .336	.490, .531
	pair rank	mean	.356, .370	<u>.423</u> , <b>.420</b>
		distr	.347, .356	.540, .548

Table 22: Pluralistic alignment error ( $\downarrow$ , Eq. 1) from MT-Bench human pairwise preferences. Comma-separated values are with and without CoT. Text styling follows Table 1. The method 'distr' uses the predicted distribution, while the other methods place probability 1 on a measure of central tendency.

Table 28 presents the intransitivity rates of different methods. Despite the capacity of QT and PS to model intransitive preferences (Savage, 1994; Finkelstein and Thorp, 2006; Conrey et al., 2013), we find that they exhibit negligible intransitivity compared to the pairwise ranking methods. Similar to Liu et al. (2024e), we observe that a stronger judge (GPT-40) exhibits less intransitivity than a weaker judge (Llama-3.1-8B). Preaggregation mean exhibits less intransitivity than post-aggregation mode. Notably, for pairwise ranking, we observe more intransitivity with CoT than without CoT, even though CoT achieves higher accuracy (Table 1).

## **G** Derivations

1526

1527

1528

1529

1530

1533

1534

1535

1537

1538

1539

1541

1542

1543

1544

1546

# G.1 Approximability of Discrete Pointwise Functions Under Discretization

**Proposition 1.** We analyze the discrete methods in Table 3. Specifically, we examine the score function r rather than  $sgn(r_1 - r_2)$ .

Let X be a random variable with support  $S \subset$ 

 $\begin{bmatrix} \frac{1}{2}, K + \frac{1}{2} \end{bmatrix} \text{ for an integer } K. \text{ Define its discretiza-}$   $\text{tion } \hat{X} \text{ by } P(\hat{X} = \hat{x}) \coloneqq P([X] = \hat{x}) \text{ for}$   $\hat{x} \in \hat{S} \coloneqq \{1, \dots, K\}, \text{ where } [\cdot] \text{ denotes round-}$  ing to the nearest integer. 1547 1548 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 1549 15

- 1. MODE may fail to be approximated: Sup-1551 pose X has a density  $f_X$  that is L-Lipschitz 1552 with  $L \leq 1$  and achieves its supremum 1553 at  $x^* \in \arg \max_{x \in S} f_X(x)$ . Let  $\hat{x}^* \in$ 1554  $\arg \max_{\hat{x}\in\hat{S}} P(\hat{X} = \hat{x}).$  Suppose  $\hat{x} \in \hat{S}$ , 1555 with arbitrarily large  $|\hat{x} - \hat{x}^*| > 1$ , satisfies 1556  $P(\hat{X} = \hat{x}^*) \ge P(\hat{X} = \hat{x}) + \frac{L}{4}$ . The above is 1557 consistent with  $[x^*] = \hat{x}$ . 1558
- 2. [MEAN] can be approximated:  $|[\mathbb{E}X] [\mathbb{E}\hat{X}]| \le 1$ .

1559

1560

1561

1562

1563

1564

1565

1567

1568

1569

1570

1571

1573

1576

1577 1578

1582

1584

- 3. MEDI and 1P can be approximated: For  $p \in (0,1), |Q_X(p) Q_{\hat{X}}(p)| \leq \frac{1}{2}$ .
- Proof.
  - 1. We present a construction.

If L = 0, the claim is immediate; assume not. Define  $d := \frac{L}{4}(\sqrt{1+8/L}-2) \ge \frac{L}{4}$ . Let  $f_X(x) = (d - \frac{L}{4}) + L(x - \hat{x} + \frac{1}{2})$  for  $x \in [\hat{x} - \frac{1}{2}, \hat{x})$ , and  $f_X(x) = (d - \frac{L}{4}) + L(\hat{x} - x + \frac{1}{2})$ for  $x \in [\hat{x}, \hat{x} + \frac{1}{2})$ , and  $f_X(x) = d + \frac{L}{4}$  for  $[x] = \hat{x}^*$ .

Around the regions  $[\hat{x} - \frac{1}{2}, \hat{x} + \frac{1}{2}), [\hat{x}^* - \frac{1}{2}, \hat{x}^* + \frac{1}{2})$ , we let  $f_X$  decrease to 0 with slope  $\pm L$ , or until reaching the domain boundary or each other. Continuity is maintained at the junction because, supposing  $\hat{x} < \hat{x}^*$  without loss of generality, the nearest endpoints  $\hat{x} + \frac{1}{2}, \hat{x}^* - \frac{1}{2}$ satisfy  $|(\hat{x} + \frac{1}{2}) - (\hat{x}^* - \frac{1}{2})| \ge 1$  and  $|f_X(\hat{x} + \frac{1}{2}) - f_X(\hat{x}^* - \frac{1}{2})| = \frac{L}{2}$ .

We verify that $P(X = \hat{x}^*) = d + \frac{L}{4} =$	1579
$P(\hat{X} = \hat{x}) + \frac{L}{4} \text{ and } \hat{x} \in {\hat{x}} \cup [\hat{x}^* - \frac{1}{2}, \hat{x}^* +$	1580
$\frac{1}{2}) = \arg\max_{x \in S} f_X(x).$	1581

It remains to check that we have a valid distribution. The total  $\int f_X$  is bounded by the case if  $f_X$  is allowed to reach 0 everywhere

Model	Method	Helpf	ulness	Corre	ctness	Cohe	rence	Comp	olexity	Verb	osity
		$W_1$	$W_2$	$W_1$	$W_2$	$W_1$	$W_2$	$W_1$	$W_2$	$W_1$	$W_2$
GPT-40	mode	.218	.311	.219	.332	.149	.252	.211	.273	.186	.257
	mean	.221	.297	.217	.318	.151	.240	.213	.262	.197	.244
	distr	.188	.279	.194	.301	.134	.233	.199	.255	.179	.249
Llama -3.1-8B	mode	.259	.369	.250	.377	.154	.280	.227	.290	.182	.255
	mean	.255	.339	.249	.347	.158	.253	.224	.274	.174	.223
	distr	.219	.328	.215	.334	.134	.250	.209	.270	.164	.234

Table 23: Pluralistic alignment error ( $\downarrow$ , Eq. 1) from HelpSteer2 human pointwise scores. No-CoT. Text styling follows Table 1. The method 'distr' uses the predicted distribution, while the other methods place probability 1 on a measure of central tendency.

Model	Method	Reward-Bench	MT-Bench
	-	.091, .105	.093, .111
	MODE	.103, .150	.128, .214
	MEAN	<u>.066</u> , .080	<u>.105</u> , .136
	[MEAN]	.104, .115	.144, .199
GPT-4o	MEDI	.101, .113	.137, .185
	1 P	.096, .117	.137, .196
	RAM	.074, .084	.111, .138
	QT	<b>.064</b> , .078	<b>.104</b> , .133
	PS	<b>.064</b> , .078	<b>.104</b> , .137
	_	.136, .063	.117, .076
	MODE	.213, .201	.223, .247
	MEAN	.149, .042	.131, <u>.048</u>
Llama -3.1-8B	[MEAN]	.213, .139	.219, .218
	MEDI	.219, .160	.224, .218
	1 P	.223, .105	.183, .133
	RAM	.168, <u>.037</u>	.156, .068
	QT	.151, <b>.034</b>	.130, <u>.048</u>
	PS	.151, .037	.129, <b>.046</b>

Table 24: Sensitivity to granularity  $(\downarrow)$  of the score distributions (Eq. 1) and of the pointwise methods computed on them (Eq. 2). Comma-separated values are with and without CoT. Text styling follows Table 1.

 $\int f_X \le P(\hat{X} = \hat{x}) + P(\hat{X} = \hat{x}^*)$ 

 $= 1 - \frac{L}{4} < 1,$ 

affecting the desired properties.

 $+\frac{1}{L}\left(d-\frac{L}{4}\right)^2+\frac{1}{L}\left(d+\frac{L}{4}\right)^2$ 

so  $f_X$  can be made a valid density by adding

an appropriately scaled uniform density, not

2. Denote the measures of  $X, \hat{X}$  as  $\mu_X, \mu_{\hat{X}}$ . The

definition of  $(X, \hat{X})$  is equivalent to the ex-

istence of a coupling  $\gamma \in \Gamma(\mu_X, \mu_{\hat{X}})$  with samples defined by  $(x, \hat{x}) \sim \gamma$  for  $x \sim \mu_X$ 

possible above:

1585

1587

1589

1592

1593

1594 1595

Model	Setting	K	MAE	MSE
GPT-4o	score	9	.090, <b>.076</b>	.057, <b>.031</b>
	score	99	.094, .095	.049, <u>.032</u>
	rank	2	.086, .087	.083, .037
	rank	3	.085, .089	.078, .035
	rank	5	.141, .182	.079, .053
Llama -3.1-8B	score	9	.199, <u>.163</u>	.125, .066
	score	99	.188, <b>.160</b>	.114, <b>.060</b>
	rank	2	.357, .329	.193, .154
	rank	3	.683, .518	.547, .340
	rank	5	.506, .342	.334, .164

Table 25: Pairwise position bias ( $\downarrow$ , see Appendix F.3) on RewardBench (see Table 26 for MT-Bench). Commaseparated values are with and without CoT. Text styling follows Table 1. We find that no-CoT always maintains or improves MSE, even when it hurts MAE.

Model	Setting	K	MAE	MSE
GPT-4o	score score	9 99	.108, <b>.091</b> .111108	.075, <b>.038</b> .066, <b>.039</b>
	rank	2	.108, .132	.100, .056
	rank rank	3 5	.108, .134 .187, .172	.093, .051 .120, .047
Llama -3.1-8B	score	9	.211, .148	.145, .056
	score	99	.193, <b>.141</b>	.129, <b>.049</b>
	rank	2	.312, .355	.174, .172
	rank	3	.618, .532	.466, .337
	rank	5	.458, .298	.293, .129

Table 26: Pairwise position bias ( $\downarrow$ ) on MT-Bench, mirroring Table 25.

and 
$$\hat{x} = [x]$$
. 159

$$\left|\mathbb{E}X - \mathbb{E}\hat{X}\right| = \left|\int (x - \hat{x}) \,\mathrm{d}\gamma(x, \hat{x})\right|$$
 159

$$\leq \int |x - \hat{x}| \,\mathrm{d}\gamma(x, \hat{x}) \leq \int \frac{1}{2} \,\mathrm{d}\gamma(x, \hat{x}) = \frac{1}{2}$$
 159

Thus, 
$$|[\mathbb{E}X] - [\mathbb{E}\hat{X}]| \le 1.$$
 1599

Space	Nectar	RM-Bench	MT-Bench
interm	<b>.086</b>	<b>.079</b>	<b>.033</b>
list	.092	.100	<u>.041</u>
direct list	.118	.105	.056

Table 27: Listwise position bias  $(\downarrow)$  with GPT-40. We report the absolute value⁵ of Spearman's  $\rho$  between the difference in the presented positions of two responses and the judgment. Text styling follows Table 1.

Model	Setting	Method	MT-Bench
GPT-4o	point score	QT	.000, .000
	point score	PS	.006, .002
	pair rank	MODE-AGG	.026, .022
	pair rank	AGG-MEAN	.007, .003
Llama -3.1-8B	point score	QT	.000, .000
	point score	PS	.001, .000
	pair rank	MODE-AGG	.234, .218
	pair rank	AGG-MEAN	.040, .023

Table 28: A study on transitivity. In each cell, we report the proportion of triplets that exhibit intransitivity, with and without CoT. (Pointwise scoring uses K = 9; pairwise ranking uses Likert-2.) In addition, our Nectar silver labels (GPT-40, Likert-5, no-CoT, mean) have an intransitivity rate of 0.020.

1601

1602

1603

1604

1605

1607

1608

1610

1611

1612

1613

1614

1615

1616

1617

3. Let  $q \coloneqq Q_X(p)$ .

$$P(\hat{X} < [q] - \frac{1}{2}) = P(X < [q] - \frac{1}{2})$$
$$\leq P(X < [q] + \frac{1}{2}) = P(\hat{X} < [q] + \frac{1}{2})$$

implying  $Q_{\hat{X}}(p) = [q]$  where  $|q - [q]| \le \frac{1}{2}$ .

< p

*Remark.* The suppositions in (1) are to impose regularity and show even then approximation may not hold. For an example of their omission, without requiring absolutely continuous X, it could place atoms at arbitrary x, preventing any margin  $P(\hat{X} = \hat{x}^*) - P(\hat{X} = \hat{x})$  less than 1 from producing an error bound. The crucial case that causes the mode to be unstable to approximate is the case of multimodality.

In (3), it is crucial that we assumed no discretization error, i.e.  $|P(\hat{X} = \hat{x}) - P([X] = \hat{x})| = 0.$ With any discretization error, we would have no bound on approximation error.

# **H** Licensing

Our usage of the artifacts below complies with their 1619 licenses. 1620

**Model Licensing** GPT-40⁶ has a proprietary li-1621 cense. Llama-3.1-8B⁷ is licensed under the Llama 1622 3.1 Community License Agreement. Mistral-1623  $7B^8$  and Prometheus-2- $7B^9$  are licensed under the 1624 Apache License 2.0.

**Dataset Licensing** The datasets contain English language data. RewardBench¹⁰ and RM-Bench¹¹ are licensed under the ODC-By license. MT-1628 Bench¹² and HelpSteer2¹³ are licensed under the CC BY 4.0 license. Nectar¹⁴ is licensed under the 1630 Apache License 2.0.

1627

1631

1632

#### Ι **Ethical Considerations**

LLMs can exhibit unwanted biases. Relying on 1633 their judgments for downstream applications can 1634 propagate these biases. Nevertheless, our findings 1635 in this paper promote practices for improving align-1636 ment with human preferences. 1637

⁸https://huggingface.co/mistralai/ Mistral-7B-Instruct-v0.3

- ¹⁰https://huggingface.co/datasets/allenai/ reward-bench
- ¹¹https://huggingface.co/datasets/THU-KEG/ **RM-Bench**
- ¹²https://huggingface.co/datasets/lmsys/mt_ bench_human_judgments

¹⁴https://huggingface.co/datasets/ berkeley-nest/Nectar

⁶https://platform.openai.com/docs/models# gpt-4o

⁷https://huggingface.co/meta-llama/Llama-3. 1-8B-Instruct

⁹https://huggingface.co/prometheus-eval/ prometheus-7b-v2.0

¹³https://huggingface.co/datasets/nvidia/ HelpSteer2/tree/main/disagreements