Improving LLM-as-a-Judge Inference with Distributional Judgment

Anonymous ACL 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 that incor-017 018 porate risk aversion can improve calibration. 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 suggest leveraging distributional output can improve 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. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

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 92 out of 120 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 can improve calibration. 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 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



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.

are widely adopted for complex tasks, improving best practices for using LLM-as-a-judge can impact many end tasks' development and evaluation.

2 Background

086

097

102

103

104

105

2.1 LLM-as-a-Judge Settings

We review three settings for LLM-as-a-judge.

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 106 the mean of the judgment distribution, mostly in 107 the pointwise setting. Liu et al. (2023b); Lee et al. (2024a); Saad-Falcon et al. (2024) note the benefits of the mean but do not empirically compare it 110 with the mode. Zawistowski (2024) shows that the 111 mean outperforms the mode for summary scoring. 112 113 Concurrent work (Yasunaga et al., 2024) shows that the mean outperforms the mode on Reward-114 Bench (Lambert et al., 2024), but the paper's focus 115 is on data-efficient alignment. Similarly, Hashemi 116 et al. (2024) show that after training a calibrator for 117

personalized alignment, the mean outperforms the mode on dialogue judgment.

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

151

152

153

155

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 text judgments. In this work, we benchmark the mode, the mean, and newly proposed methods for leveraging distributional outputs 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). In this work, we analyze the effect of CoT in tandem with the 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.

Model	odel Setting Met		RewardBench				MT-Bench	
	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		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 4a	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, <u>87.4</u>	<u>86.2, 86.5</u>
OF 1-40	pair score	mean	<u>97.2, 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>
L lama-3 1-8B	nair 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-5.1-0D	pair score	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
	point score	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	nair score	mode	87.6, <u>89.9</u>	<u>40.2</u> , <u>40.4</u>	<u>74.0, 73.0</u>	67.4, 72.4	67.3, 68.9	<u>79.3, 79.8</u>
Wilstrai-7D	pan 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>
	nair rank	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 rank	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>
	pan 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>51.0</b> , <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 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 92 out of 120 cases. No-CoT outperforms CoT in 30 out of 40 cases when using the mean for pointwise or pairwise scoring.

# 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

156

157

158

159

160

161

162

163

164

165

166

167

168

To infer a judgment from the LLM judge'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

170ModelsAs LLM judges, we use gpt-4o-2024-08-17106 (shortened to GPT-40) (OpenAI et al., 2024),172Llama-3.1-8B-Instruct (Llama-3.1-8B) (Dubey173et al., 2024), Mistral-7B-Instruct-v0.3 (Mistral-7B)174(Jiang et al., 2023), and Prometheus-2-7B (Kim175et al., 2024). We cover a commonly used closed-176sourced LLM (GPT-40), as well as smaller open-177sourced 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.

178

179

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

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.

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.

Model	Setting	RewardBench	MT-Bench
GPT-4o	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.

# 3.3 Results

203

210

211

212

213

214

215

216

217

218

221

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 92 out of 120 cases. 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.

**CoT often harms LLM-as-a-judge** For the scoring settings, no-CoT outperforms CoT in 30 out of 40 cases when using the mean. For the pairwise ranking setting, CoT outperforms no-CoT, except with GPT-40 on RewardBench.

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). 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. We report the standard deviation of judgment distributions in Table 2, which confirms this trend.

Which setting works the best? Comparing different LLMs, we find GPT-40 performs better with 229 pairwise judgment (e.g. 88.0 for pointwise scoring 230 vs. 90.5 for pairwise ranking on RewardBench) as in prior work, but the smaller models often do better with pointwise judgment and rely heavily on 234 CoT for pairwise ranking (e.g. with Prometheus-2-7B on MT-Bench, 75.4 $\rightarrow$ 33.4 when removing CoT 235 from pairwise ranking, compared to 81.6 with no-CoT pointwise scoring). We believe this is because pairwise judgment demands a more powerful judge 238

to leverage the context. Thus, in pairwise ranking with the smaller models, the reasoning gained by CoT often outweighs the distributional signal lost in the process. Nonetheless, using pairwise scoring (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.

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

262

263

264

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

283

284

287

# 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.

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

Additional Metric: Mean Squared Error For this further study, 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% (Section 3.2). 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 briefly motivate each newly introduced method below.

We consider the rounded mean as a discrete variant of the mean for a discreteness-controlled comparison with the mode.

1P and RAM reflect risk aversion. 1P takes an approach contrary to MODE; instead of focusing on

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). Discrete/continuous and more properties are explained in Section 4.1.

Model	Method	Rewar	dBench	MT-	MT-Bench	
	metalou	$Acc \uparrow$	$MSE\downarrow$	$Acc \uparrow$	$MSE\downarrow$	
	MODE	83.5	.115	80.0	.137	
	MEAN	88.0	.102	<u>83.2</u>	.097	
	[MEAN]	85.2	.109	80.2	.146	
CPT 4a	MEDI	84.6	.112	80.2	.142	
OF 1-40	1 P	84.3	.116	81.0	.138	
	RAM	88.4	.087	83.4	.104	
	QT	87.9	.096	83.2	.118	
	PS	87.8	.096	<u>83.3</u>	.103	
	MODE	72.2	.191	71.8	.143	
	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	.129	<u>81.4</u>	.109	
	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 10 for CoT). Text styling follows Table 1.

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  $\sigma_{-}$ .

PS has the intuitive interpretation of 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.

#### 4.2 Results

Main Results Table 4 shows that the top pointwise methods are the continuous ones (MEAN, RAM, QT, PS), in both accuracy and MSE. Even among discrete methods, MODE has particularly low accuracy and without substantially different MSE. RAM outperforms MEAN on RewardBench (e.g. 0.087 vs. 0.102 MSE with GPT-40), suggest-

Model	Method	Tie	rate	MEAN's accuracy		
		K = 9	K = 99	K = 9	K = 99	
	MODE	.20	.21	71	73	
CDT 4a	[MEAN]	.16	.03	67	53	
OF 1-40	MEDI	.17	.09	70	62	
	1 P	.16	.08	66	60	
	MODE	.35	.24	69	70	
Llama -3.1-8B	[MEAN]	.26	.07	64	61	
	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 11 for CoT and Table 12 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% indicates room for improving accuracy, and excess of 75% indicates room for improving MSE.

ing that risk aversion can be helpful for preference modeling.

Study: Score Granularity and Ties We show here that ties explain the finding above that the discrete methods fall beyond 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 11), indicating that the performance gap is well explained by ties. Tie behavior varies by method; in particular, MODE has the most ties and the highest MEAN accuracy, amounting to the most untapped signal for determining the better response.

Table 5 further shows that granularizing the

305

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

322



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).

	Method	Rewar	RewardBench		Bench
		Acc $\uparrow$	$MSE\downarrow$	Acc $\uparrow$	$MSE\downarrow$
	MODE	81.5-2.0	.132+.017	78.1-1.9	.152+.015
	MEAN	<b>86.7</b> _{-1.3}	$.108_{+.006}$	$82.9_{-0.3}$	<b>.099</b> _{+.002}
0	[MEAN]	$86.5_{+1.3}$	$.127_{+.018}$	$82.7_{+2.5}$	$.182_{+.036}$
7	MEDI	$85.2_{\pm 0.6}$	$.126_{+.014}$	81.5+1.3	$.170_{+.028}$
Ļ,	1 P	$86.4_{+2.1}$	$.116_{+.000}$	$82.7_{+1.7}$	.165+.027
0	RAM	<b>86.7</b> _{-1.7}	<b>.089</b> +.002	<b>83.0</b> _{-0.4}	$.105_{+.001}$
	QT	<u>86.6</u> _{-1.3}	$.114_{+.018}$	$82.7_{-0.5}$	$.147_{+.029}$
	PS	<u>86.6</u> _{-1.2}	$.105_{+.009}$	82.4-0.9	$.107_{+.004}$
	MODE	71.9-0.3	.221+.030	75.1+3.3	$.168_{+.025}$
В	MEAN	$79.3_{+0.0}$	$.156_{+.001}$	$81.3_{-0.2}$	<b>.103</b> ₀₀₁
-8	[MEAN]	$78.5_{+3.5}$	$.198_{+.012}$	80.7+5.7	$.180_{+.035}$
ς.	MEDI	$76.5_{+2.9}$	$.207_{+.016}$	$80.1_{+6.2}$	$.161_{+.019}$
Llama-	1 P	$78.5_{+2.5}$	$.195_{+.012}$	$81.5_{+2.3}$	$.177_{+.030}$
	RAM	<b>79.7</b> _{-0.2}	<b>.129</b> _{+.000}	$81.1_{-0.3}$	$.109_{+.000}$
	QT	$78.7_{-0.3}$	$.177_{+.013}$	81.3+0.2	$.143_{+.027}$
	PS	78.6-0.3	$.163_{+.002}$	<b>81.8</b> +0.4	$.111_{+.001}$

Table 6: Pointwise results over methods (K = 99). No-CoT (see Table 13 for CoT). Subscripts denote change from K = 9 (Table 4). Text styling follows Table 1.

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.

323

324

325

326

327

330

334

335

336

339

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 also suffers the highest MSE among discrete methods on RewardBench.

Taken together, Tables 4, 5, 6 show that even in use cases where discrete scores are desired, one should consider alternatives to the mode. **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.

340

341

343

344

345

350

351

352

354

355

356

357

360

361

362

363

365

367

369

372

#### 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).

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

	Method	Reward	dBench	MT-Bench		
		Acc $\uparrow$	$MSE\downarrow$	Acc $\uparrow$	$MSE\downarrow$	
	mode agg	88.1, 89.3	.099, .090	<u>86.1</u> , 84.9	.139, .142	
0	agg mode	88.4, <u>90.3</u>	.112, .094	<u>86.5</u> , 85.2	.154, .154	
7	medi agg	88.1, 89.3	.099, .091	<u>86.1,</u> 84.9	.138, .142	
Ę.	agg medi	88.4, 90.0	.111, .094	<b>86.6</b> , 85.4	.153, .146	
0	mean agg	88.9, <b>90.4</b>	.098, <b>.077</b>	<u>86.5</u> , 85.4	.132, .100	
	agg mean	88.9, <b>90.4</b>	.098, <u>.078</u>	<b>86.6</b> , 85.4	.132, <b>.097</b>	
В	mode agg	56.7, 52.4	.240, .279	57.5, 53.4	.192, .176	
-8	agg mode	73.1, 66.1	.265, .337	<b>78.1</b> , 70.9	.222, .268	
	medi agg	56.8, 52.5	.240, .279	57.5, 53.5	.192, .176	
Llama-	agg medi	72.9, 65.3	.261, .319	<u>78.0</u> , 69.1	.218, .238	
	mean agg	<u>73.2</u> , 65.6	<b>.207</b> , .229	<b>78.2</b> , 70.5	<b>.144</b> , <u>.146</u>	
	agg mean	<b>73.2</b> , 66.3	.222, .240	<u>78.1</u> , 70.8	.155, .155	

Table 7: Pairwise ranking results over methods, using Likert-3. We denote pre- or post-aggregation by prepending or appending 'agg', respectively. Commaseparated 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, .087	
Llomo	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 8: Pairwise ranking results over Likert-K scales, using pre-aggregation mean. Comma-separated values are with and without CoT. Text styling follows Table 1.

#### 5.2 Methods Results

373

374

375

382

384

390

391

393

Table 7 shows that, especially with Llama-3.1-8B, accuracy depends little on the measure of central tendency and mostly on when we aggregate, with aggregating first leading to higher accuracy (as much as 56.7 vs. 73.1 for post- vs. pre-aggregation using the mode with Llama-3.1-8B on Reward-Bench). Considering that the timing of aggregation does not affect accuracy if the two runs agree, this shows that even for inconsistent judgments caused by position bias, there is 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 at-

tribute 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. Furthermore, with GPT-40, we find that although CoT has differing effects on accuracy for RewardBench and MT-Bench, the minimum MSE for either is achieved with no-CoT, highlighting the discord between CoT's sharpening effect and calibration.

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

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 the simplest one, Likert-2, performs 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-5.¹ Even so, 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).

Similar to Table 7, we find that in every case with GPT-40, no-CoT achieves lower MSE than CoT, even in cases where CoT leads to higher accuracy. This result is in line with AlpacaEval, which uses no-CoT and judgment probabilities, but deviating from WB-Reward and Arena-Hard-Auto, which use CoT and decoded judgments.

# 6 Listwise Judgment

Listwise judgment is not as prevalent as pointwise or pairwise judgment, but is more efficient (Zhu et al., 2024). Furthermore, listwise prompting grants the judge the maximal context for comparison (Buyl et al., 2023), generalizing the advantage of pairwise prompting over pointwise prompting, given a sufficiently capable judge.

¹We remark that the evaluation setups in these works and ours differ in that theirs use coarse-grained, model-level agreement using Chatbot Arena Elo scores (Zheng et al., 2023b), while ours uses fine-grained, instance-level agreement.

#### 6.1 Judgment Spaces and Methods

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

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 pairwise² preferences from one of the following judgment spaces using the mode (textual output) or the mean (distributional output).

- 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.

**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, i.e. 9 pairwise preferences. MT-Bench contains 160 prompts,

Space	Method	Nectar		RM-Bench		MT-Bench	
		Acc	MSE	Acc	MSE	Acc	MSE
interm	mode	80.4	.155	62.1	.339	80.8	.201
	mean	80.4	.048	62.5	.243	<u>80.7</u>	.121
list	mode	82.2	.156	62.4	.376	83.7	.189
	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.

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, 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 (Section 3.3), 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.

# 7 Conclusion

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. 501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

483

 $^{^{2}}$ We retain the pairwise evaluation setup from previous sections; see Appendix D.1 for discussion.

# 521 Limitations

530

531

532

533

534

535

546

547

549

551

552

553

554

560

561

562

563 564

565

567

568

522Downstream PerformanceIn this paper, we523evaluate LLM-as-a-judge design decisions by their524performance on preference modeling datasets.525However, this setup may not reveal downstream526impacts. We do not explore the impact of distribu-527tional judgments on reinforcement learning from528AI feedback (RLAIF) (Lee et al., 2024a) or human529decision 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.

**CoT** We conclude from our results that CoT often hurts judgment performance. However, we only consider one prompt design per setting for eliciting CoT reasoning (Appendix C) and do not perform prompt optimization. Furthermore, we do not consider more extensive test-time scaling, such as asking the judge to produce its own reference response (Zheng et al., 2023b) or aggregating many CoT judgment runs (Zhang et al., 2024a; Stureborg et 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,

Nicholas Joseph, Saurav Kadavath, John Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv*, abs/2204.05862.

569

570

571

572

573

574

575

576

577

578

579

580

582

584

585

586

587

588

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

- 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.
- 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, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien

623 Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, 641 Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, 648 Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona 657 Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier 659 Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, 667 Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, 670 Shaoliang Nie, Sharan Narang, Sharath Raparthy, 671 Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gu-673 674 rurangan, Sydney Borodinsky, Tamar Herman, Tara 675 Fowler, Tarek Sheasha, Thomas Georgiou, Thomas 676 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong 677 Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor 678 Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 679 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue 685 Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh,

Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre

687

690

691

692

693

694

695

696

697

698

699

702

705

707

708

709

712

713

714

715

716

717

718

719

720

721

722

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

Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.

751

752

754

771

773

774

776

778

779

790

797

- 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*.
- 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.

809

810

811

812

813

814

815

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

- 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, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 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*.
- 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.

966

967

968

969 970

971

915

- 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.

874

878

897

901

902

903

904

905

906

907

908

909

910 911

912

913

914

- 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.
- 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*.
- 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.
- 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, and S. Longpre. 2023. Octopack: Instruction tuning code large language models. *ArXiv*, abs/2308.07124.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson,

972 Chak Li, Chan Jun Shern, Channing Conger, Char-973 lotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris 974 Koch, Christian Gibson, Christina Kim, Christine 976 Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-983 dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, 993 Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya 997 Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub 999 Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, 1000 Jason Kwon, Jason Phang, Jason Teplitz, Jason 1002 Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-1003 avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui 1004 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Lan-1005 1006 ders, Joel Parish, Johannes Heidecke, John Schul-1007 man, Jonathan Lachman, Jonathan McKay, Jonathan 1008 Uesato, Jonathan Ward, Jong Wook Kim, Joost 1009 Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, 1010 Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai 1011 Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin 1012 Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, 1013 1014 Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, 1015 Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle 1016 Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-1017 ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia 1018 Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-1019 ian Weng, Lindsay McCallum, Lindsey Held, Long 1020 Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, 1021 1022 Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine 1023 Boyd, Madeleine Thompson, Marat Dukhan, Mark 1024 Chen, Mark Gray, Mark Hudnall, Marvin Zhang, 1025 Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, 1026 Max Johnson, Maya Shetty, Mayank Gupta, Meghan 1027 Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao 1028 Zhong, Mia Glaese, Mianna Chen, Michael Jan-1029 ner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, 1030 Miguel Castro, Miguel Oom Temudo de Castro, 1031 Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-1032 nal Khan, Mira Murati, Mo Bavarian, Molly Lin, 1033 1034 Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-1035 talie Cone, Natalie Staudacher, Natalie Summers,

Natan LaFontaine, Neil Chowdhury, Nick Ryder, 1036 Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, 1037 Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, 1039 Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, 1040 Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Pe-1043 ter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip 1044 Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming 1045 Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra-1046 jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul 1047 Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob 1049 Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-1050 dani, Romain Huet, Rory Carmichael, Rowan Zellers, 1051 Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agar-1054 wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu 1056 Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi-1057 rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, 1061 Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas 1063 Shadwell, Tianhao Zheng, Todd Underwood, Todor 1064 Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce 1066 Walters, Tyna Eloundou, Valerie Oi, Veit Moeller, 1067 Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne 1068 Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, 1069 Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, 1070 Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen 1071 He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and 1072 Yury Malkov. 2024. Gpt-40 system card. Preprint, 1073 arXiv:2410.21276. 1074

Vishakh Padmakumar, Chuanyang Jin, Hannah Rose Kirk, and He He. 2024. Beyond the binary: Capturing diverse preferences with reward regularization. *ArXiv*, abs/2412.03822.

1075

1076

1077

1078

1079

1081

1083

1084

1085

1086

1087

1088

1089

1090

1091

- 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 1094

1095 1096	as preference machines for few-shot ordinal classifi- cation. <i>ArXiv</i> , abs/2408.03359.	helps mainly on math and symbolic reasoning. <i>ArXiv</i> , abs/2409.12183.	1149 1150
1097	Vyas Raina, Adian Liusie, and Mark Gales. 2024. Is	Rickard Stureborg, Dimitris Alikaniotis, and Yoshi	1151
1098	llm-as-a-judge robust? investigating universal adver-	Suhara. 2024. Large language models are inconsis-	1152
1099	sarial attacks on zero-shot llm assessment. arXiv	tent and biased evaluators. ArXiv, abs/2405.01724.	1153
1100	preprint arXiv:2402.14016.	Sijun Tan Siyuan Zhuang Kyle Montgomery	115/
4404	Deventh Canai Daddy, Jaallysaal, Daa, Vifai Vu	William V Tang Alejandro Cuadron Chenguang	1154
1101	Md Arefet Sulten Decure Swein Avimum Sil and	Wang Raluca Ada Pona and Ion Stoica 2024	1156
1102	Heng Ji 2024 First: Faster improved listwise	Iudgebench: A benchmark for evaluating llm-based	1157
1103	regarking with single token decoding arXiv preprint	judges, arXiv preprint arXiv:2410.12784.	1158
1105	arXiv:2406.15657		
1105	u/Mtv.2+00.15057.	Raphael Tang, Xinyu Crystina Zhang, Xueguang Ma,	1159
1106	Paul Röttger, Hannah Rose Kirk, Bertie Vidgen,	Jimmy Lin, and Ferhan Ture. 2023. Found in the mid-	1160
1107	Giuseppe Attanasio, Federico Bianchi, and Dirk	dle: Permutation self-consistency improves listwise	1161
1108	Hovy. 2023. Xstest: A test suite for identifying exag-	ranking in large language models. In North Amer-	1162
1109	gerated safety behaviours in large language models.	ican Chapter of the Association for Computational	1163
1110	ArXiv, abs/2308.01263.	Linguistics.	1164
		T N Tideman 1987 Independence of clones as a	1165
1111	Jon Saad-Falcon, Rajan Vivek, William Berrios, Nan-	criterion for voting rules Social Choice and Welfare	1166
1112	dita Shankar Naik, Matija Franklin, Bertie Vidgen,	4(3):185–206	1167
1113	Amanpreet Singh, Douwe Kiela, and Shikib Mehri.	(5).105 200.	
1114	2024. Lmunit: Fine-grained evaluation with natural	Xuezhi Wang and Denny Zhou. 2024. Chain-of-	1168
1115	language unit tests. Arxiv, abs/2412.15091.	thought reasoning without prompting. ArXiv,	1169
1116	Swarnadeen Saha, Omer Levy, Asli Celikvilmaz, Mohit	abs/2402.10200.	1170
1117	Bansal Jason Weston and Xian Li 2023 Branch-		
1118	solve-merge improves large language model evalua-	Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov,	1171
1119	tion and generation. ArXiv, abs/2310.15123.	and Ilmoiny Baldwin. 2023. Do-not-answer: A	11/2
		abs/2208 12287	1173
1120	Richard P. Savage. 1994. The paradox of nontransitive	aus/2508.15587.	1174
1121	dice. American Mathematical Monthly, 101:429-	Zhilin Wang, Alexander Bukharin, Olivier Delal-	1175
1122	436.	leau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Olek-	1176
		sii Kuchaiev, and Yi Dong. 2024a. Helpsteer2-	1177
1123	Markus Schulze. 2011. A new monotonic, clone-	preference: Complementing ratings with preferences.	1178
1124	independent, reversal symmetric, and condorcet-	ArXiv, abs/2410.01257.	1179
1125	consistent single-winner election method. Social		
1126	<i>Choice and Welfare</i> , 36(2):267–303.	Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi	1180
1197	Ammar Shaikh Rai Abhiiit Dandekar Sreedath Panat	Zeng, Gerald Shen, Daniel Egert, Jimmy Zhang,	1181
1128	and Rai Abhijit Dandekar, 2024. Cheval: A frame-	Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev.	1182
1129	work for evaluating and interpreting cognitive biases	2024b. Helpsteer2: Open-source dataset for	1183
1130	in llms. <i>ArXiv</i> , abs/2412.03605.	training top-performing reward models. ArXiv,	1184
	····,····	aus/2400.08075.	COLL
1131	Lin Shi, Weicheng Ma, and Soroush Vosoughi. 2024.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	1186
1132	Judging the judges: A systematic investigation of	Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou.	1187
1133	position bias in pairwise comparative assessments by	2022. Chain of thought prompting elicits reasoning	1188
1134	llms. arXiv preprint arXiv:2406.07791.	in large language models. ArXiv, abs/2201.11903.	1189
4405	Around Siththomorium Consider I ville and D 1		
1135	Anand Situtharanjan, Cassidy Laidiaw, and Dylan	Michihiro Yasunaga, Leonid Shamis, Chunting Zhou,	1190
1136	Hadileid-Menell. 2023. Distributional preference	Andrew Cohen, Jason Weston, Luke Zettlemoyer,	1191
1100	context in rlbf, arYiu preprint arYiu:2312.08358	and Marjan Ghazvininejad. 2024. Alma: Alignment	1192
1130	context in fini. <i>urxiv preprint urxiv.2512.06556</i> .	with minimal annotation.	1193
1139	Taylor Sorensen. Jared Moore, Jillian Fisher	Chen Ye. Wei Xiong Yuheng Zhang Nan Jiang and	110/
1140	Mitchell L Gordon, Niloofar Mireshghallah, Christo-	Tong Zhang 2024a Online iterative reinforcement	1195
1141	pher Michael Rytting, Andre Ye, Liwei Jiang,	learning from human feedback with general prefer-	1196
1142	Ximing Lu, Nouha Dziri, et al. 2024. Position: A	ence model.	1197
1143	roadmap to pluralistic alignment. In Forty-first		
1144	International Conference on Machine Learning.	Seonghyeon Ye, Doyoung Kim, Sungdong Kim,	1198
		Hyeonbin Hwang, Seungone Kim, Yongrae Jo,	1199
1145	Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez,	James Thorne, Juho Kim, and Minjoon Seo. 2023.	1200
1146	Dongwei Jiang, Manya Wadhwa, Prasann Singhal,	Flask: Fine-grained language model evaluation	1201
1147	Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Dur-	based on alignment skill sets. arXiv preprint	1202
1148	rett. 2024. To cot or not to cot? chain-of-thought	arXiv:2307.10928.	1203
	1	4	
	1		

- 1204 1205 1206 1207 1208 1209 1210 1211 1212 1213 1214 1215 1216 1217 1218
- 1217 1218 1219 1220 1221
- 1221 1222 1223
- 1224
- 1225 1226 1227 1228
- 1229 1230 1231 1232 1233 1234 1235 1236
- 1230 1237 1238 1239 1240 1241
- 1242 1243 1244 1245
- 1246 1247
- 1248 1249
- 1250 1251

1252 1253 1254

1255 1256

1257 1258

- 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*.
- 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, Joseph E. Gonzalez, Ion Stoica, and Haotong Zhang. 2023a. Lmsys-chat-Im: 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,<br/>and G. Zuccon. 2023. A setwise approach for effec-<br/>tive and highly efficient zero-shot ranking with large<br/>language models. In Annual International ACM SI-<br/>GIR Conference on Research and Development in<br/>Information Retrieval.1259<br/>1260Information Retrieval.1261<br/>1262

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1287

1288

1289

1290

1291

1292

1293

1295

1296

1297

1299

1300

1301

1302

# A Related Work

# A.1 LLM-as-a-judge 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] 1304 that signifies the advantage of  $A_1$  over  $A_2$ . For 1305 accuracy, we take the sign of the prediction. For 1306 MSE, we rescale predictions from [-1, 1] to [0, 1]. 1307

1343 1344

1345

1346

1348

1349

1352

1353

1354

1356

1340 1341

1328

1321 1322

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 prompts for the various settings are in Appendix C.

# **B.1** Pointwise Methods

The pointwise methods are 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 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  $\sigma$  term grants continuity. Specifically, let  $k, k' \in \{1, \ldots, K\}$  with  $k \neq k'$ . Let  $X_1$  have a two-point distribution  $(1 - \epsilon)\delta_k + \epsilon \delta_{k'}$  and let  $X_2$  have a deterministic distribution  $\delta_k$ , for  $\epsilon \in$ [0,1]. Then as a function of  $\epsilon$ , MEAN $(X_1, X_2)$  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, OT 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**

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  $\mathbf{o} \in \mathbf{O} \coloneqq \{(1,2), (2,1)\}$ , 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  $P(X_{o_1}^{\mathbf{o}}, X_{o_2}^{\mathbf{o}})$ , we could compute the marginals and use any method in Table

3. However, the judge first outputs the score for 1357  $A_{o_1}$  and conditions on it when outputting the score 1358 for  $A_{o_2}$ , i.e.  $X_{o_1}^{\mathbf{o}}$  and  $X_{o_2}^{\mathbf{o}}$  are not independent. 1359 Thus, the full joint distribution  $P(x_{o_1}, x_{o_2}) =$ 1360  $P(x_{o_1})P(x_{o_2} \mid x_{o_1})$  can only be obtained by injecting each  $x_{o_1} \in \{1, \ldots, K\}$  into the context to 1362 access  $P(x_{o_2} \mid x_{o_1})$ . This is feasible with local 1363 models but not with API-access models where in-1364 ference cost scales with K. Hence, we stick to a single run and condition on the greedily decoded 1366  $x_{o_1} = \arg \max_k P(X_{o_1}^{\mathbf{o}} = k)$ , giving us 1367

$$X_{\Delta}^{\mathbf{o}} \stackrel{d}{=} (X_1^{\mathbf{o}} - X_2^{\mathbf{o}}) \mid (X_{o_1}^{\mathbf{o}} = x_{o_1})$$
 1368

as a proxy for the score difference  $X_1^{\mathbf{o}} - X_2^{\mathbf{o}}$ . Semantically,  $X^{\mathbf{o}}_{\Delta}$  is symmetric (i.e. there should be 1370 no default preference for  $A_1$  or  $A_2$ ), so we would 1371 like our scalar judgment to be some measure of *cen*-1372 tral tendency (mode, median, or mean). As shown in Figure 2, we also have the choice of whether 1374 to aggregate the judgments from the two orders of 1375 presentation before or after computing the measure 1376 of central tendency. 1377

For pre-aggregation, we simply take the mixture distribution,

$$P(X_{\Delta} = \delta) \coloneqq \frac{1}{|\mathbf{O}|} \sum_{\mathbf{o} \in \mathbf{O}} P(X_{\Delta}^{\mathbf{o}} = \delta)$$
 1380

1378

1379

1387

1388

1389

1390

for all  $\delta \in \{-(K-1), \ldots, K-1\}$ , leaving more 1381 sophisticated approaches such as the convolution and Wasserstein barycenter for future study: 1383

$$AGG-MODE \coloneqq \operatorname{sgn}(\operatorname{mode}(X_{\Delta}))$$
 1384

$$AGG-MEDI \coloneqq \operatorname{sgn}(\operatorname{median}(X_{\Delta}))$$
 1385

$$\operatorname{AGG-MEAN} \coloneqq \operatorname{MEAN}(X_{\Delta}),$$
 1386

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:

$$\text{MODE-AGG} \coloneqq \frac{\sum_{\mathbf{o} \in \mathbf{O}} \text{mode}(X_{\Delta}^{\mathbf{o}})}{\sum_{\mathbf{o} \in \mathbf{O}} |\text{mode}(X_{\Delta}^{\mathbf{o}})|}$$
139

$$\text{MEDI-AGG} \coloneqq \frac{\sum_{\mathbf{o} \in \mathbf{O}} \operatorname{median}(X_{\Delta}^{\mathbf{o}})}{\sum_{\mathbf{o} \in \mathbf{O}} |\operatorname{median}(X_{\Delta}^{\mathbf{o}})|}$$
1392

$$MEAN-AGG := \frac{1}{|\mathbf{O}|} \sum_{\mathbf{o} \in \mathbf{O}} MEAN(X_{\Delta}^{\mathbf{o}}),$$
 1393

taking  $\frac{0}{0} \coloneqq 0$ . 1394

# B.2.2 Pairwise Ranking

We prompt the LLM judge to express its prefer-1396 ence on a K-point Likert scale: [>, <] (Likert-2), 1397 [>, =, <] (Likert-3), or  $[\gg, >, =, <, \ll]$  (Likert-1398 5). Assigning the symbols  $[\gg, >, =, <, \ll]$  the 1399 numerical values [2, 1, 0, -1, -2], the methods for 1400 pairwise ranking then follow those above for pair-1401 wise scoring. We remark that the 'mode' and 1402 'mean' for pairwise scoring and pairwise rank-1403 ing in Table 1 are with post-aggregation and pre-1404 1405 aggregation, respectively.

# 1406 B.3 Listwise Methods

1407 The listwise methods are introduced in Section 6.1.

# 1408 C Prompts

1409We present representative example prompts to illus-<br/>trate the different settings. The prompts are adapted1410from MT-Bench (Zheng et al., 2023b). Auxiliary<br/>modifications are not shown, such as the prompt1413for second-turn evaluation in MT-Bench.

# 1414 C.1 Judgment Extraction Details

To identify the token position containing the judg-1415 ment, we use the specified format when available 1416 (e.g. "Rating A: {rating_a}." in pairwise scor-1417 ing). Otherwise, we use the latest token position 1418 with more than 0.5 total probability assigned to 1419 judgment tokens. If no valid token is found, we de-1420 fault the judgment to the minimum score of 1 in the 1421 scoring setting, and to a tie in the ranking setting. 1422 (For Nectar experiments, we exclude instances with 1423 1424 invalid silver-label judgments.)

1425For the local models (Llama-3.1-8B, Mistral-7B,1426Prometheus-2-7B) in no-CoT prompting, we force1427a prefix of the assistant's response (e.g. "Rating1428A: ") and use a single output token as the judgment1429token position.

# C.2 Pointwise Prompts

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

Please act as an impartial judge and

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

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

Please act as an impartial judge and

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

1432

1433

# 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]

# C.3 Pairwise Prompts

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

Please act as an impartial judge and

- $\ \ \, \hookrightarrow \ \ \, evaluate \ the \ quality \ of \ the \ responses$
- $\rightarrow$  provided by two AI assistants to the
- $\hookrightarrow \quad$  user prompt displayed below. Your
- $\, \hookrightarrow \quad evaluation \ should \ consider \ factors$
- $\label{eq:such as the helpfulness, relevance,} \quad such as the helpfulness, relevance,$
- $\hookrightarrow$  accuracy, depth, creativity, level of
- $\, \hookrightarrow \,$  detail, and ethicality of their responses.
- $\, \hookrightarrow \ \ \, Begin \ your \ evaluation \ by \ comparing$
- $\, \hookrightarrow \, \ \ \, \text{the two responses and provide a short} \,$
- $\hookrightarrow$  explanation. Avoid any position biases
- $\rightarrow$  and ensure that the order in which the
- $\, \hookrightarrow \quad responses \ were \ presented \ does \ not$
- $\, \hookrightarrow \quad \text{influence your decision. Do not allow}$
- $\hookrightarrow \quad \text{the length of the responses to} \quad$
- $\ \ \, \hookrightarrow \quad influence \ your \ evaluation. \ Do \ not$
- $\, \hookrightarrow \,$  favor certain names of the assistants.
- $\hookrightarrow$  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
- $\rightarrow$  integer scores from 1 to 9.

1435

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.

# System prompt for pairwise ranking (Likert-5, CoT)

Please act as an impartial judge and

- $\label{eq:evaluate} \to \quad \text{evaluate the quality of the responses}$
- $\hookrightarrow$  provided by two AI assistants to the
- $\, \hookrightarrow \,$  user prompt displayed below. You
- $\hookrightarrow$  should choose the assistant that
- $\hookrightarrow$  follows the user's instructions and
- $\hookrightarrow$  answers the user's question better.
- $\hookrightarrow$  Your evaluation should consider
- $\hookrightarrow$  factors such as the helpfulness,
- $\label{eq:constraint} \hookrightarrow \quad \text{relevance, accuracy, depth, creativity,}$
- $\, \hookrightarrow \,$  level of detail, and ethicality of their
- $\hookrightarrow$  responses. Begin your evaluation by
- $\hookrightarrow$  comparing the two responses and
- $\, \hookrightarrow \,$  provide a short explanation. Avoid any
- $\rightarrow$  position biases and ensure that the
- $\, \hookrightarrow \,$  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
- $\hookrightarrow$  of the assistants. Be as objective as
- $\rightarrow$  possible. After providing your
- $\hookrightarrow$  explanation, output your final verdict
- $\rightarrow$  by strictly following this format:
- $\hookrightarrow$  "[[>>]]" if assistant A is significantly
- $\rightarrow$  better, "[[>]]" if assistant A is slightly
- $\rightarrow$  better, "[[=]]" for a tie, "[[<]]" if
- $\hookrightarrow$  assistant B is slightly better, and
- $\rightarrow$  "[[<<]]" if assistant B is significantly
- $\rightarrow$  better.

# 1441

# 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

1444The listwise prompts are adapted from Nectar (Zhu1445et al., 2024).

# System prompt for listwise judgment (N = 7), with 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
- $\, \hookrightarrow \ \ \, quality \ of \ the \ completions \ provided \ by$
- $\rightarrow$  the 7 AI assistants. Your evaluation
- $\, \hookrightarrow \,$  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
- $\leftrightarrow$  the order in which to perform pairwise
- $\hookrightarrow$  comparisons. Output an array in
- $\leftrightarrow$  which, for each pairwise comparison,
- $\hookrightarrow$  you choose the letter of the better
- $\rightarrow$  response, or '=' for a tie. The array
- $\hookrightarrow$  should be comma-separated and
- $\rightarrow$  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: [[_, _,
- → _, _, _, _, _]] 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
- $\hookrightarrow$  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
- $\rightarrow$  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
- $\hookrightarrow$  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),
- $\rightarrow$  strictly in the following format: [[_, _,
- $\rightarrow$  _, _, _, _, _]] where '_' contains an
- $\hookrightarrow$  assistant's letter name.

Avoid any position biases and ensure that

- $\hookrightarrow$  the order in which the responses were
- $\rightarrow$  presented does not influence your
- $\hookrightarrow$  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.

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]

1448

D

Datasets

1449

1450

1451

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

safety. Each instance consists of a prompt, a cho-1452 sen response, and a rejected response, all manually 1453 verified. The dataset categories are Chat, with 358 1454 instances sourced from AlpacaEval (Li et al., 2023) 1455 and MT-Bench (Zheng et al., 2023b); Chat Hard, 1456 with 456 instances sourced from MT-Bench and 1457 LLMBar (Zeng et al., 2023); Safety, with 740 in-1458 stances sourced from XSTest (Röttger et al., 2023), 1459 Do-Not-Answer (Wang et al., 2023), and original 1460 data; and Reasoning, with 1431 instances sourced 1461 from PRM800k (Lightman et al., 2023) and Hu-1462 manEvalPack (Muennighoff et al., 2023). Except 1463 for excluding the prior sets category, we follow the 1464 original work and compute the final score as the 1465 average of the category scores. 1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

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

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519 1520

1521

1522

1523

1524

1526

1527

1528

1529

1531

1532

1533

1534

1535

1536

1537

1539

1540

1541

1542

1543

1545 1546 For the listwise setting, we use the same evaluation setup as with the pointwise and pairwise setting.³ 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.

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

Tables 10 (K = 9) and 13 (K = 99) show pointwise results over methods (expanded versions of Tables 4 and 6). Tables 11 and 12 show expanded tie analyses on RewardBench (simplified in Table 5) and MT-Bench.

# F Analysis

#### 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  $0.1/1.1 \approx 0.091$ . 1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1564

1565

1566

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1596

Table 14 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 disagreement is correlated with the uncertainty in the LLM's judgment distribution. On datasets with multiple human judgments per instance, we compute Spearman's  $\rho$  between the standard deviation of the human judgments and that of the LLM's judgment distribution.

For MT-Bench, we take the 961 instances with multiple human judgments. Table 15 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.

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 16 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.

# F.1.3 Pluralistic Alignment

We finally evaluate the alignment between predicted judgment distributions and human judgment

³This means that the MT-Bench results can be directly compared across the settings.

Model	Model Method		dBench	MT-E	MT-Bench	
		Acc $\uparrow$	$MSE\downarrow$	Acc $\uparrow$	$MSE\downarrow$	
GPT-4o	MODE MEAN [MEAN] MEDI 1P RAM QT PS	85.1, 83.5 87.4, 88.0 85.1, 85.2 85.0, 84.6 84.8, 84.3 87.4, <b>88.4</b> 87.4, 87.9 87.4, 87.8	.116, .115 .099, .102 .116, .109 .116, .112 .120, .116 <b>.072</b> , .087 .107, .096 .106, .096	82.0, 80.0 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>	.150, .137 .115, <b>.097</b> .150, .146 .150, .142 .141, .138 .114, .104 .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, .191 .198, .155 .238, .186 .238, .191 .238, .183 .146, <b>.129</b> .220, .164 .216, .161	74.8, 71.8 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, .143 .129, <b>.104</b> .172, .145 .176, .142 .172, .147 .126, .109 .154, .116 .149, .110	

Table 10: 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, .20	.09, .21	64, 71	61, 73	+0.0, -0.1	-0.0, -0.1
CDT 4a	[MEAN]	.13, .16	.02, .03	65, 67	61, 53	+0.0, +0.0	+0.0, -0.0
UF 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	+0.1, -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 11: 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% indicates room for improving accuracy, and excess of 75% indicates room for improving MSE. 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.

distributions. We quantify the distance between two distributions with the Wasserstein *p*-distance for  $p \in \{1, 2\}$  (Eq. 1). A higher *p* more heavily punishes large point distances |x - y|. We scale the judgment spaces to [0, 1] so that  $W_p(\mu, \nu) \in [0, 1]$ .

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

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

Table 17 shows that using a distributional prediction has little success in improving alignment with the MT-Bench human pairwise preferences, but Table 18 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. 1614

1615

1616

1617

1618

1619

#### F.2 Sensitivity to Score Granularity

Adopting the view that LLMs latently encode a con-1620 tinuous distribution but output a discretization of 1621 it (Gillman et al., 2024), we analyze how faithfully 1622 functions of the (latent) continuous distribution can 1623 be approximated by those functions computed on 1624 the (observed) discretization. For practical inter-1625 est, this manifests as robustness to the choice of 1626 K, with convergence in distribution to the continu-1627 ous distribution as  $K \to \infty$ . Thus, independently of the "principledness" of certain functions of a 1629

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, .24	.09, .24	61, 64	62, 67	+0.0, +0.2	+0.1, -0.6
CDT 4a	[MEAN]	.13, .21	.02, .04	61,64	42, 48	+0.0, +0.0	+0.0, +0.0
OF 1-40	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	66, 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 12: Tie analysis for discrete pointwise methods on MT-Bench, mirroring Table 11.

Model Method		Rewa	rdBench	MT-Bench		
		Acc ↑	MSE ↓	Acc ↑	MSE ↓	
	MODE	86.0+0.9, 81.5-2.0	.117 _{+.001} , .132 _{+.017}	$83.9_{+1.9}, 78.1_{-1.9}$	.147003, .152+.015	
	MEAN	87.4 _{+0.0} , 86.7 _{-1.3}	$.097_{002}, .108_{+.006}$	$84.8_{+1.2}, 82.9_{-0.3}$	.105010, .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	1P	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}	<b>.071</b> ₀₀₁ , .089 _{+.002}	<b>85.1</b> _{+1.2} , 83.0 _{-0.4}	$.108_{006}, .105_{+.001}$	
	QT	<u>87.3</u> _{-0.1} , <u>86.6</u> _{-1.3}	$.112_{+.005}, .114_{+.018}$	$84.8_{+1.3}, 82.7_{-0.5}$	$.149_{+.010}, .147_{+.029}$	
	PS	$\underline{87.3}_{-0.1}, \underline{86.6}_{-1.2}$	$.105_{001}, .105_{+.009}$	$\underline{84.8}_{+1.3}, 82.4_{-0.9}$	$.130_{006}, .107_{+.004}$	
	MODE	73.4+3.8, 71.9-0.3	.222015, .221+.030	77.3+2.5, 75.1+3.3	.190+.013, .168+.025	
	MEAN	$75.9_{+3.2}, \underline{79.3}_{+0.0}$	$.183_{015}, .156_{+.001}$	$79.3_{+0.6}, \underline{81.3}_{-0.2}$	.125 ₀₀₄ , <b>.103</b> ₀₀₁	
	[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}	<u>.132</u> 014, <b>.129</b> +.000	$79.7_{+0.9}, 81.1_{-0.3}$	.123003, .109+.000	
	QT	$75.7_{+2.9}, 78.7_{-0.3}$	$.214_{006}, .177_{+.013}$	$78.8_{\pm 0.1}, 81.3_{\pm 0.2}$	$.179_{+.025}, .143_{+.027}$	
	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 13: 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 10). Text styling follows Table 1.

Model	Setting	K	RewardBench	MT-Bench
	point score	9	.000, .008	.000, .012
CDT 4a	point score	99	.362, .409	.357, .440
GP 1-40	pair rank	3	.000, .018	.000, .019
	pair rank	5	.014, .049	.021, .041
	point score	9	.009, .040	.013, .025
Llama -3.1-8B	point score	99	.356, .379	.382, .365
	pair rank	3	.044, .091	.051, .081
	pair rank	5	.107, .194	.107, .245

Table 14: A study on multimodality (see Appendix F.1.1). Comma-separated values are with and without CoT.

ground-truth continuous distribution, it is appropriate to examine the effect of discretization on our ability to approximate them to begin with. Our theoretical result is stated in Proposition 1 (see Appendix G.1 for full statement, proof, and discussion).

1630

1631

1632

1633

1634

1635

1636

1638

**Proposition 1.** Among the discrete methods in Table 3, MODE computed on continuous distributions may fail to be approximated by the same function

Model Setting MT-Bench +0.21, +0.24 point score GPT-40 +0.19, +0.27pair score pair rank +0.19, +0.27 point score +0.21, +0.14Llama pair score +0.20, +0.24 -3.1-8B +0.02, -0.04pair rank

Table 15: Spearman's $\rho$ between standard deviation of
human judgments and that of LLM's judgment distri-
bution. 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 $[P(X_1 > X_2), P(X_1 =$
$X_2$ , $P(X_1 < X_2)$ ].

computed on their discretizations, even under regularity conditions. Meanwhile, [MEAN], MEDI, and 1P admit an approximation error bound.

We empirically assess the robustness to K of1642the score distributions produced by the LLM judge1643as well as the functions computed on them. The1644former is not addressed by Proposition 1, which1645

1639

1640

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 16: 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
		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>
		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 17: 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.

assumes the score distributions to be errorless discretizations and thus consistent across granularities.

# F.2.1 Sensitivity of Score Distributions

1646

1647

1648

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1663

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

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

where  $\Gamma(\mu, \nu)$  is the set of couplings of  $\mu$  and  $\nu$ .

## 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

$$FR \coloneqq \frac{\|\operatorname{sgn}(\mathbf{a}^9) - \operatorname{sgn}(\mathbf{a}^{99})\|_1}{\|\operatorname{sgn}(\mathbf{a}^9)\|_1 + \|\operatorname{sgn}(\mathbf{a}^{99})\|_1} \in [0, 1].$$
(2) 1666

1665

1668

1669

1670

1671

1672

1673

1674

1675

1676

1677

1678

1679

1680

1681

1682

1684

1685

1686

1687

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

## F.2.3 Results

Table 19 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.

**Evaluation Metrics** For the pairwise setting (scoring or ranking), we measure mean absolute error (MAE) and mean squared error (MSE) between the two judgments from the two orders, using pre-aggregation mean. Compared to MAE, MSE punishes a few large errors more than many small errors.

For the listwise setting, we measure Spearman's  $\rho$  between the difference in the presented positions of two responses and the judgment.

⁴In every judgment space, GPT-40 tends to favor responses

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

Table 18: 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-40	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, .048
Llomo	[MEAN]	.213, .139	.219, .218
-3 1-8B	MEDI	.219, .160	.224, .218
-5.1-6D	1 P	.223, .105	.183, .133
	RAM	.168, .037	.156, .068
	OT	.151, .034	.130, .048
	PS	.151, .037	.129, .046

Table 19: 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.

**Results** Tables 20 and 21 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 22 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.

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

Model	Setting	K	MAE	MSE
	score	9	.090, <b>.076</b>	.057, <b>.031</b>
	score	99	.094, .095	.049, <u>.032</u>
GPT-40	rank	2	.086, .087	.083, .037
	rank	3	.085, .089	.078, .035
	rank	5	.141, .182	.079, .053
	score	9	.199, <u>.163</u>	.125, .066
Llama	score	99	.188, <b>.160</b>	.114, <b>.060</b>
-3.1-8B	rank	2	.357, .329	.193, .154
	rank	3	.683, .518	.547, .340
	rank	5	.506, .342	.334, .164

Table 20: Pairwise position bias ( $\downarrow$ , see Appendix F.3) on RewardBench (see Table 21 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
	score	9	.108, <b>.091</b>	.075, <b>.038</b>
	score	99	.111, .108	.066, <b>.039</b>
GPT-40	rank	2	.108, .132	.100, .056
	rank	3	.108, .134	.093, .051
	rank	5	.187, .172	.120, .047
	score	9	.211, .148	.145, .056
Llama	score	99	.193, <b>.141</b>	.129, <b>.049</b>
-3.1-8B	rank	2	.312, .355	.174, .172
	rank	3	.618, .532	.466, .337
	rank	5	.458, .298	.293, .129

Table 21: Pairwise position bias  $(\downarrow)$  on MT-Bench, mirroring Table 20.

#### 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

1722

1723

1724

1725

1726

1727

1728

1729

that are presented earlier.

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

Table 22: Listwise position bias ( $\downarrow$ ) with GPT-40. We report the absolute value⁴ 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 point score pair rank pair rank	QT PS MODE-AGG AGG-MEAN	.000, .000 .006, .002 .026, .022 .007, .003
Llama -3.1-8B	point score point score pair rank pair rank	QT PS MODE-AGG AGG-MEAN	.000, .000 .001, .000 .234, .218 .040, .023

Table 23: 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.

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.

Table 23 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 than1758post-aggregation mode. Notably, for pairwise rank-1759ing, we observe more intransitivity with CoT than1760without CoT, even though CoT achieves higher ac-1761curacy (Table 1).1762

# **G** Derivations

# G.1 Approximability of Discrete Pointwise Functions Under Discretization

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

Let X be a random variable with support  $S \subset [\frac{1}{2}, K + \frac{1}{2})$  for an integer K. Define its discretization  $\hat{X}$  by  $P(\hat{X} = \hat{x}) := P([X] = \hat{x})$  for  $\hat{x} \in \hat{S} := \{1, \dots, K\}$ , where  $[\cdot]$  denotes rounding to the nearest integer.

- 1. MODE may fail to be approximated: Suppose X has a density  $f_X$  that is L-Lipschitz with  $L \leq 1$  and achieves its supremum at  $x^* \in \arg \max_{x \in \hat{S}} f_X(x)$ . Let  $\hat{x}^* \in \arg \max_{\hat{x} \in \hat{S}} P(\hat{X} = \hat{x})$ . Suppose  $\hat{x} \in \hat{S}$ , with arbitrarily large  $|\hat{x} \hat{x}^*| > 1$ , satisfies  $P(\hat{X} = \hat{x}^*) \geq P(\hat{X} = \hat{x}) + \frac{L}{4}$ . The above is consistent with  $[x^*] = \hat{x}$ .
- 2. [MEAN] can be approximated:  $|[\mathbb{E}X] [\mathbb{E}\hat{X}]| \le 1$ .
- 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  $P(\hat{X} = \hat{x}^*) = d + \frac{L}{4} =$  1802  $P(\hat{X} = \hat{x}) + \frac{L}{4} \text{ and } \hat{x} \in \{\hat{x}\} \cup [\hat{x}^* - \frac{1}{2}, \hat{x}^* +$  1803  $\frac{1}{2}) = \arg \max_{x \in S} f_X(x).$  1804 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 possible above:

1809 
$$\int f_X \le P(\hat{X} = \hat{x}) + P(\hat{X} = \hat{x}^*)$$

1805

1806

1807

1808

1812

1813

1814

1815

1817

1818

1819

so  $f_X$  can be made a valid density by adding an appropriately scaled uniform density, not affecting the desired properties.

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

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

Denote the measures of X, X̂ as μ_X, μ_{X̂}. The definition of (X, X̂) is equivalent to the existence of a coupling γ ∈ Γ(μ_X, μ_{X̂}) with samples defined by (x, x̂) ~ γ for x ~ μ_X and x̂ = [x].

1820 
$$|\mathbb{E}X - \mathbb{E}\hat{X}| = \left| \int (x - \hat{x}) \, \mathrm{d}\gamma(x, \hat{x}) \right|$$
$$\leq \int |x - \hat{x}| \, \mathrm{d}\gamma(x, \hat{x}) \leq \int \frac{1}{2} \, \mathrm{d}\gamma(x, \hat{x}) = \frac{1}{2}$$

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

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

1822

1823

1825

1826

182

1828 1829

1830 1831

1832

1833

1834

1835

1836

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

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

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

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.

1837In (3), it is crucial that we assumed no discretiza-1838tion error, i.e.  $|P(\hat{X} = \hat{x}) - P([X] = \hat{x})| = 0.$ 1839With any discretization error, we would have no1840bound on approximation error.

# H Licensing

Our usage of the artifacts below complies with their1842licenses.1843

1841

1844

1845

1846

1847

1850

1851

1852

1853

1854

1855

**Model Licensing** GPT- $40^5$  has a proprietary license. Llama- $3.1-8B^6$  is licensed under the Llama 3.1 Community License Agreement. Mistral- $7B^7$  and Prometheus- $2-7B^8$  are licensed under the Apache License 2.0.

**Dataset Licensing** The datasets contain English language data. RewardBench⁹ and RM-Bench¹⁰ are licensed under the ODC-By license. MT-Bench¹¹ and HelpSteer2¹² are licensed under the CC BY 4.0 license. Nectar¹³ is licensed under the Apache License 2.0.

# I Ethical Considerations

LLMs can exhibit unwanted biases. Relying on 1856 their judgments for downstream applications can 1857 propagate these biases. Nevertheless, our findings 1858 in this paper promote practices for improving alignment with human preferences. 1860

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

¹³https://huggingface.co/datasets/ berkeley-nest/Nectar

⁵https://platform.openai.com/docs/models#
gpt-4o

⁶https://huggingface.co/meta-llama/Llama-3. 1-8B-Instruct

⁸https://huggingface.co/prometheus-eval/ prometheus-7b-v2.0

⁹https://huggingface.co/datasets/allenai/ reward-bench

¹⁰https://huggingface.co/datasets/THU-KEG/ RM-Bench

¹¹https://huggingface.co/datasets/lmsys/mt_ bench_human_judgments

¹²https://huggingface.co/datasets/nvidia/ HelpSteer2/tree/main/disagreements