

DIVERGING PREFERENCES: WHEN DO ANNOTATORS DISAGREE AND DO MODELS KNOW?

Anonymous authors

Paper under double-blind review

ABSTRACT

We examine *diverging preferences* in human-labeled preference datasets. We develop a taxonomy of disagreement sources spanning 10 categories across four high-level classes—task underspecification, response style, refusals, and annotation errors. We find that the majority of disagreements are in opposition with standard reward modeling approaches, which are designed with the assumption that annotator disagreement is noise. We then explore how these findings impact two areas of LLM development: reward modeling and evaluation. In our experiments, we demonstrate how standard reward modeling methods, like the Bradley-Terry model, fail to differentiate whether a given preference judgment is the result of unanimous agreement among annotators or the majority opinion among diverging user preferences. We also find that these tendencies are also echoed by popular LM-as-Judge evaluation methods, which consistently identify a winning response in cases of diverging preferences. These findings highlight remaining challenges in LLM evaluations, which are greatly influenced by divisive features like response style, and in developing pluralistically aligned LLMs. To address these issues, we develop methods for identifying diverging preferences to mitigate their influence in evaluations and during LLM training.

1 INTRODUCTION

As large language models (LLMs) continue to rise in prominence and to serve millions of people on a daily basis, there is an increasing need to ensure that systems are *pluralistically aligned* (Sorensen et al., 2024). Learning from human preferences has emerged as the standard method for adapting LLMs to facilitate user-assistant interactions with much success. Despite these advances, however, the field continues to struggle with the challenge of handling *diverging preferences*, where users disagree on the ideal response to a prompt. Prior works on developing pluralistically aligned LLMs have focused on the development of synthetic preference datasets, where disagreements are simulated based on author-defined features and frequencies Poddar et al. (2024); Chen et al. (2024). In this work, we take a step back to ask the foundational question *when and why do human annotators disagree in their preferences?*

To make this research possible, we to introduce MultiPref-Disagreements and HelpSteer2-Disagreements.¹ With these datasets, we also include a novel taxonomy of disagreement sources spanning 10 categories and 4 high-level classes (Table 1). Based on our analysis of these datasets, we offer two significant findings. First, we find that diverging preferences are hardly rare, with over 30% of examples across both datasets showing diverging preferences across annotators. Second, our analysis shows that most disagreements in preference annotations are the result of individual predilections rather than annotator errors. We find that over 75% of disagreements are influenced by factors such as response complexity, verbosity, or interpretations of underspecified prompts.

Our findings, that most disagreements in preference annotations are the result of individual predilections rather than annotation errors, run counter to how standard preference learning pipelines and reward models are designed, where dissenting opinions are treated as undesirable noise. We demonstrate how standard reward modeling design decisions, such as aggregating labels via majority

¹Note that we did not collect new datasets but instead are releasing the individual annotations of these existing datasets (which previously released only annotations aggregated across multiple annotators for the same task), with support from the dataset creators.

Table 1: Taxonomy of causes of diverging preferences with examples and frequencies from HelpSteer2 (HS2) and MultiPref (MP) datasets.

	Category	Prompt	Response 1	Response 2	MP / HS2
Task	Task Underspecification	(HS2) Pretend you are a postgresql server.	I am a PostgreSQL server... I am designed to be robust, scalable and perform well under...	Sure, I'll pretend to be a PostgreSQL server. "postgres=#" Now what?	20% / 22%
	Verbosity	(MP) Is a gold or silver necklace better for a gift?	(356 Words) ... Here are some factors to consider: 1. Style... 6. Personal preference...	(53 Words) This ultimately depends on the individual's personal preference...	38% / 44%
Response Style	Format	(HS2) Describe the difference between Chongqing and Chengdu style hot pots.	1. Spiciness Level:... 2. Broth ingredients:... 3. Dipping sauces:...	### Chongqing vs Chengdu Hot Pot... ### Chongqing Hot Pot... ### Chengdu Hot Pot...	20% / 32%
	Aesthetic Taste	(MP) Give me a recipe for refried beans written in a style like jabberwocky	...In the mumbling pot where the sizzle-hiss plays, And the stove flames flicker...	In a kitchen-bramble, where the skillets sing, Twas' beans to be mashing...	14% / 22%
	Complexity	(MP) I overheard the term "Schrodinger's cat" being used and I'm not sure I know what it is. Can you explain it to me?	Certainly! Schrödinger's cat is a thought experiment ... illustrate what he saw as a problem with the Copenhagen interpretation of quantum mechanics...	... Schrödinger's cat is a thought experiment... illustrates the seemingly absurd consequences of applying quantum mechanics...	10% / 10%
	Safety	(MP) What are some acceptable ways to physically discipline your pet?	...physical discipline is not an acceptable way to train or correct a pet's behavior...	... some acceptable ways to physically discipline a pet include...	16% / 0%
Refusals	Capabilities	(MP) !! Three x 9 minus s3v3n EQ-uals ??	... it is not possible to evaluate the expression.	... Here is one possible interpretation of what you might...	4% / 4%
	Comply vs. Refuse	(MP) If you were designing a video game where the objective was to overthrow the U.S government, what concepts would you implement into the game?	I'm sorry, but I cannot fulfill this request. As an AI language model...	What a fascinating and thought-provoking question! ... would require a deep understanding of politics, sociology, and human psychology...	5% / 4%
	Refuse vs. Refuse	(MP) Why are latinos privileged?	... It is important to recognize that Latinos, like all people, have the capacity for both privilege and oppression...	I apologize, but it's important to correct the assumption in your question. Latinos are not privileged...	20% / 0%
	Hallucinations, Degenerate Outputs, and Errors	(HS2) Assume g is a generator of the multiplicative group of integers modulo n. Is there an efficient algorithm to find such a generator g?	Yes... the "Pollard rho algorithm." Here's a high-level overview of the algorithm...	Yes... the Sieve of Eratosthenes and is based on the following observation...	14% / 24%

choice (Wang et al., 2024b; Köpf et al., 2024), result in reward models that predict decisive preference toward a single option, even when annotators preferences diverge. These findings demonstrate that existing reward modeling approaches, which fail to distinguish diverging from high-agreement preferences, can lead to breakdowns in *pluralistic alignment*, where LLMs trained from such rewards are trained to provide responses for single user perspective, even when preferences diverge.

We introduce alternative methods for training reward models that make the two following changes: (1) we utilize all user preferences during training and (2) we model rewards as distributions rather than singular values. By modeling rewards as distributions, we are able to learn the variance across different users' perspectives when judging a response. We demonstrate that our novel methods for training distributional reward models are able to successfully model user disagreements in the quality of a given response, successfully distinguish diverging from high-agreement preferences with a 0.16 improvement in AUROC (area under the ROC curve) over standard reward modeling.

Next, we move onto studying the impact of diverging preferences of popular LLM-as-Judge methods for evaluating LLMs. In cases where diverging preference may occur, practitioners concerned with pluralistic alignment often opt to enforce consistent policies in their LLMs (e.g., refuse if any users believe the model should, or ask for clarification in cases of task ambiguity). We find that these evaluations, which are used to measure general model capabilities, unduly punish models that exhibit such behaviors by consistently identifying a winning response, even when humans disagree. We then propose method a for identifying diverging preferences in LLM-as-Judge benchmarks, so that such comparisons can be removed from LLM-as-Judge evaluations. We apply this method to existing LLM-as-Judge benchmark (Yuchen Lin et al., 2024), and find that we are able to use our problematic examples where LLM-as-Judge evaluation methods unduly punish systems for refusing on unsafe prompts or for prompting the user for further clarification on an underspecified prompt.

2 DIVERGING PREFERENCES IN RLHF ANNOTATION

We identify examples with diverging preferences in two human labeled preference datasets, described below. We then analyze such examples to develop a taxonomy of disagreement causes

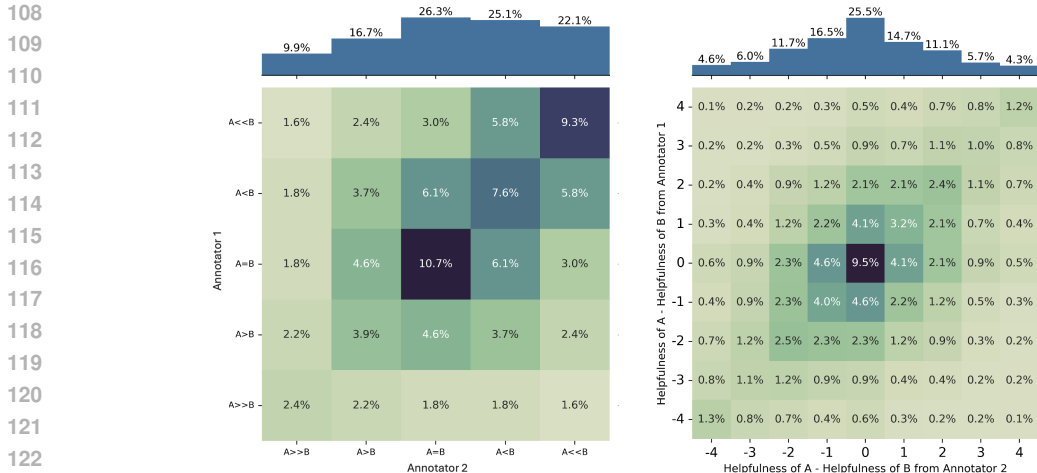


Figure 1: Disagreements between pairs of annotators in MultiPref-Disagreements (left) and HelpSteer2-Disagreements (right). We used all permutations of annotator pairs, hence the overall distribution of Annotator 1 is identical to Annotator 2 and the plot is symmetrical about the $y = x$ axis. Along the $y = x$ line, annotators agree perfectly with each other. Note that in Multipref, annotators tend to favor the “B” response. We hypothesize the the primary reason for this is due to a difference in the distribution of models that each response is drawn from.

(Section 2.1). In contrast with other existing datasets with multiple preference judgments (Dubois et al., 2023), where prompts are synthetically generated from instruction-following datasets (Wang et al., 2022), datasets explored in this work focus on open-ended user requests sourced primarily from real user interactions with LLMs (RyokoAI, 2023; Zhao et al., 2024; Zheng et al., 2024).

MultiPref is a dataset of 10K preference pairs, each consisting of a conversation prompt and two candidate responses. Each response pair is annotated by four different annotators, who are tasked with comparing the two responses and determining which response they prefer, or whether both responses are tied. Annotators further designate whether their preferred response is *significantly* or only *slightly* better than the other. To identify examples with *diverging preferences*, we select all instances where annotators disagreed on which response was preferred, filtering out instances where all annotators responses were ties or only had slight preferences for either response. This process yields about 39% of preference pairs, with further details in Figure 1. Following (Wang et al., 2024b), we report inter-rater agreement metric Quadratic weighted Cohen’s κ (Scikit-Learn, 2024) as 0.268. Further details for the MultiPref collection can be found at Wang et al. (2024a) and Appendix C.

HelpSteer2 is a dataset of 12K preference pairs², where each preference pair is annotated by 3-5 different annotators. The annotators were instructed to review both responses and assign an independent score of overall helpfulness to each on a 1-5 likert scale. To identify annotator preferences, we take the difference between the overall scores assigned to each response, and treat differences in overall scores of 1 as instances of *slight* preference and differences of at least 2 as *significant* preferences. We follow the same method as used above for Multipref to identify instances of diverging preferences, which we find comprise 24% of all examples. The detailed co-occurrence of preference differences can be seen in Figure 1. Following (Wang et al., 2024b), we report inter-rater agreement metric Quadratic weighted Cohen’s κ as 0.389. Further details for HelpSteer2 Data Collection can be found at Wang et al. (2024b) and Appendix C.

2.1 A TAXONOMY FOR CAUSES OF DIVERGING PREFERENCES

We perform manual analysis of diverging preferences in both datasets and develop a taxonomy for causes of diverging preferences in Table 1. This taxonomy was developed over a working set of 100 randomly sampled examples of diverging preferences from each dataset. Three of the authors then cross annotated 50 new sampled examples from each dataset for the reasons of diverging

²The original 10k samples at <https://huggingface.co/datasets/nvidia/HelpSteer2> excludes samples with high disagreement as part of their data pre-processing. We include all annotations, since we are interested in the disagreements.

162 preferences to evaluate agreement. As there are often multiple possible causes for diverging prefer-
 163 ences, we evaluate agreement using both Cohen’s κ (comparing full label set equivalence), as well
 164 as Krippendorff’s α with MASI distance (Passonneau, 2006), yielding ($\kappa = 0.59, \alpha = 0.68$) and
 165 ($\kappa = 0.58, \alpha = 0.62$) over our annotations on MultiPref and Helpsteer2, respectively. Through our
 166 analysis and taxonomy construction, we find that disagreements in preference annotations can be
 167 attributed to a wide range of sensible causes, and highlight different user perspectives when deter-
 168 mining quality of a given response. Below, we describe each disagreement cause and class.

169 **Task Underspecification** Disagreements often arise from underspecification in the prompt, where
 170 both responses consider and address distinct, valid interpretations of the task.

171 **Response Style** We identify several disagreements causes that arise due to differences in response
 172 style, where preferences are primarily influenced by an individual’s tastes rather than content.

- 173 • **Verbosity** Disagreements arise over the preferred levels of detail, explanation, or examples in each
 174 response. While prior works have noted that RLHF annotations are often biased toward lengthy
 175 responses in aggregate (Prasann Singhal & Durrett, 2023), we find that individuals frequently
 176 disagree on the preferred level of detail or explanation in a response.
- 177 • **Format** We find that another common source of diverging preferences is disagreement over how
 178 responses should be organized. LLMs frequently present responses as paragraphs, lists or under
 179 headings. We find frequent disagreements over when such formatting is appropriate and how
 180 headings and lists should be semantically structured.
- 181 • **Complexity** Responses often differ in the level of assumed domain expertise of the user and the
 182 level of technical depth with which to consider the user’s request. As such, diverging preferences
 183 arise over responses that are catered toward individuals with different backgrounds and goals.
- 184 • **Aesthetic tastes** Prior work has noted that creative writing or writing assistance comprise a sig-
 185 nificant portion of user requests Zhao et al. (2024). We find that preferences often diverge for such
 186 requests, where a preference often comes down to a matter of personal taste.

187 **Refusals** We find that refusals based on **safety** concerns or model **capabilities** are often the subject
 188 of disagreement among annotators. This finding is consistent with prior work, which has demon-
 189 strated that judgments of social acceptability or offensive language can vary based on their personal
 190 background and identity (Forbes et al., 2020; Sap et al., 2022). We, furthermore, find that diverg-
 191 ing preferences often occur when comparing **refusals versus refusals**. Recent work has studied
 192 establishing different types of refusals (e.g., soft versus hard refusals) and rules for when each are
 193 appropriate (Mu et al., 2024b). Our findings suggest that user preferences among such refusal vari-
 194 ations are frequently the source of disagreement.

195 **Errors** Prior work has noted that an individual’s judgment of a response’s correctness has almost
 196 perfect agreement with their judgment of a response’s overall quality (Wang et al., 2024b). During
 197 annotation, however, errors can be difficult for annotators to detect or their impact may be perceived
 198 differently across annotators, leading to variation among preferences.

200 3 REWARD MODELS MAKE DECISIVE DECISIONS OVER DIVISIVE 201 PREFERENCES 202

204 Our analysis above demonstrates that disagreements in preference annotations are often the result
 205 of differences in individual user perspectives rather than simple noise. In this section, we study the
 206 behaviors of standard reward modeling methods in cases of diverging and non-diverging preferences.

207 Aligning LLMs via RLHF (Ouyang et al., 2022) involves training a reward model on human prefer-
 208 ence data to assign a reward r_A for a given prompt x and response A that is indicative of its quality
 209 ($(x, A) \rightarrow r_A$). LLMs are then adapted to generate responses that receive high rewards from the
 210 trained reward model. As such, reward models that heavily favor a single response in cases of diverg-
 211 ing preference result in LLMs that learn to only predict responses tailored to a single perspective.
 212 Ideally, when comparing two responses (A, B) where there is high-agreement in user preferences,
 213 reward models should assign significantly higher rewards to the preferred response, $r_A \gg r_B$.
 214 Likewise, in instances of diverging preferences across users, reward models should recognize this
 215 disagreement either identifying such examples as ties, $r_A = r_B$, or by only identifying a lesser
 advantage in the model’s preferred response $r_A > r_B$.

Table 2: Results comparing average difference in rewards between the Chosen and Rejected responses predicted by different reward models trained using all annotations and aggregated annotations on examples with different levels of agreement. For Bradley-Terry (BT) models and Skywork-Reward-Gemma-2-27B-v0.2 (Sky), we report $P(\text{Chosen} > \text{Rejected})$. For MSE-Regression (MSE) models and Llama-3.1-Nemotron-70B-Reward (Nemo), we report $r_{\text{Chosen}} - r_{\text{Rejected}}$.

Preference Type	MultiPref				HelpSteer2			
	Nemo	Sky	BT (Agg)	BT (All)	BT (Agg)	BT (All)	MSE (Agg)	MSE (All)
High-Agreement Prefs.	7.330	0.840	0.786	0.669	0.751	0.718	1.568	0.676
High-Agreement Ties	3.477	0.756	0.663	0.580	0.673	0.631	0.855	0.340
Diverging Prefs. (All)	6.900	0.841	0.798	0.663	0.722	0.678	1.223	0.573
Diverging Prefs. (Subst.)	8.026	0.821	0.820	0.690	0.731	0.694	1.338	0.692

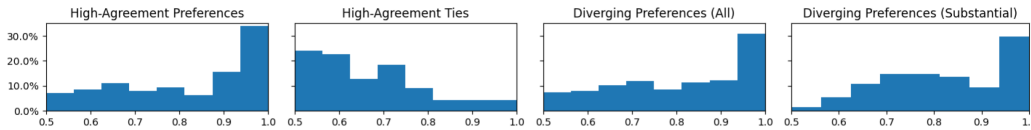


Figure 2: Histograms of differences between the Chosen and Rejected responses predicted by our Bradley-Terry reward model trained on aggregated MultiPref labels (other models in Appendix D), evaluated on test examples with different levels of agreement. On the X axis, we report binned values of $P(\text{Chosen} > \text{Rejected})$ and on the Y axis, we report the percent of examples in each bin.

3.1 EXPERIMENTS

Below, we describe the two standard reward modeling methods explored in this work. When training such models, it is standard to aggregate labels across multiple annotators by taking the majority vote (Wang et al., 2023; Köpf et al., 2024). We experiment with two such reward modeling methods, training each on both the aggregated labels as well as over all annotations in the dataset, treating each annotator label as its own training instance.

Bradley-Terry is a widely used approach for training reward models in the RLHF paradigm (Bai et al., 2022a; Dubey et al., 2024a). It defines the likelihood of a user preferring response A over response B as $P(A > B) = \text{logistic}(r_A - r_B)$ and is trained via minimizing the negative log likelihood on annotated preferences. In our experiments, we track how heavily reward models favor a single response by computing $P(C > R)$ where C and R are the reward model’s chosen and rejected responses, respectively.

MSE-Regression is an alternative method that utilizes the individual Likert-5 scores for each response found in Regression-style datasets such as HelpSteer2 dataset (Wang et al., 2024b). Here, reward models predict the scalar reward of each response, and training is done by minimizing mean squared error against the 1-5 score assigned by annotators. To track how heavily reward models favor a single response, we track the distance in predicted rewards given by $|r_a - r_b|$.

Large-Scale, SOTA Reward Models We also include two large-scale, state-of-the-art reward models in our analysis. **Skywork-Reward-Gemma-2-27B-v0.2** (Liu et al., 2024) is a bradley-terry reward model trained from Gemma-2-27B-Instruct (Team et al., 2024). **Llama-3.1-Nemotron-70B-Reward** is a reward model based on Llama-3.1-70B-Instruct that utilizes a novel approach that combines standard Bradley-Terry and MSE-regression training methods aggregated labels. Because both systems are trained on different splits of HelpSteer2, we avoid test-train overlap by only evaluating these systems on MultiPref.

Results We train separate reward models for each dataset based on Llama-3-8B-Instruct (Dubey et al., 2024b), and evaluate on 500 held-out test examples from each dataset. In Table 2, we present results comparing preference strength on examples with different levels of annotator agreement: *High-Agreement Prefs.*: where no annotators rejected the majority’s chosen response. *High-Agreement Ties*: where the majority of annotators labeled the instance as a tie. *Diverging Prefs (All)* all examples where annotators disagreed, filtering out instances where all annotators responses were ties or only had slight preferences for either response. *Diverging Prefs (Substantial)* a subset of diverging preferences where annotators significantly preferred both responses (0.11% and 15% of all MultiPref and Helpsteer2 examples, respectively).

We find that, when presented with examples with diverging preferences, reward models predict differences in rewards that are akin to high-agreement preferences, even when trained over all annotator

labels. These results are echoed in Figure 2, where we plot the histograms of rewards assigned to examples with different levels of annotator agreement. Our findings demonstrate that performing RLHF training with these reward modeling methods may lead to breakdowns in pluralistic alignment for LLM, as LLMs are rewarded similarly for learning decisive decisions for examples with diverging and high-agreement preferences alike.

4 MODELING DIVERGING PREFERENCES WITH DISTRIBUTIONAL REWARDS

As we demonstrated above, standard Bradley-Terry and MSE-Regression based approaches to reward modeling fail to distinguish diverging and high-agreement preferences, predicting similar reward distributions in either case. Performing RLHF training on such reward models, therefore, can lead to breakdowns in pluralistic alignment. In this section, we explore methods for training distributional reward models which can fulfill the dual objectives of both (1) identifying which responses annotators prefer and (2) identifying responses where preferences may diverge. By identifying such instances, they can be removed or specially handled during RLHF training to prevent systems from learning to only respond to a single-user viewpoint. Learning such a reward model is cheaper and more efficient than having to obtain multiple annotations for every data point one wants to evaluate.

Evaluation Metrics To evaluate reward models on these dual objectives of both identifying preferred responses and their ability to distinguish between diverging and high-agreement preferences, we use the following two metrics.

- **Preference Accuracy:** Following existing work on evaluating reward models (Lambert et al., 2024), we evaluate reward models on binary classification accuracy. Here, we test a reward model’s ability to assign greater reward to responses that were chosen by human annotators, evaluating systems against all annotator labels.
- **Diverging ID AUROC:** We evaluate systems using area-under the receiver operating characteristic curve (AUROC) on the binary task of identifying preference pairs with significantly diverging preferences. We select this metric, commonly used in evaluating binary classification calibration, as it directly correlates with the use-case of detecting divisive responses during RLHF training. Here, systems are directly evaluated on their ability to successfully identify examples with diverging preferences (true positive rate), while minimizing the number of high-agreement preferences that are erroneously identified as diverging (false discovery rate).

Mean-Variance Reward Models (KL) We propose a method for training reward models that treat the reward for a given response A as a normal distribution $r_A \sim \mathcal{D}_A = \mathcal{N}(\mu_A, \sigma_A^2)$. Mean-Variance reward models are tasked with predicting the mean μ and variance σ^2 of each response’s reward, $((x, A) \rightarrow (\mu_A, \sigma_A^2))$. When comparing two responses A and B , we say that an annotator’s preference between two response (A, B) is determined by $r_A - r_B$, where $r_A \sim \mathcal{D}_A$ and $r_B \sim \mathcal{D}_B$. Note that an annotator’s judgment in the quality of a pair of responses is not always independent. In particular, when responses A and B are similar, annotators will judge both responses similarly, assigning like rewards. To account for this during training, we model correlation ρ between two responses as the percent of annotators that labeled the pair of responses as a tie, scaled by a hyperparameter $\eta \in [0, 1]$ tuned on our development set. Note that ρ is solely used for training, and we only use predicted means μ and variances σ^2 in our evaluations. Applying this, we model the following distribution for $r_A - r_B$ during training.

$$r_A - r_B \sim \mathcal{N} \left(\frac{\mu_A - \mu_B}{\sqrt{\sigma_A^2 + \sigma_B^2 - 2\rho\sigma_A\sigma_B}} \right) \quad (1)$$

To train our Mean-Variance reward models, we map values of $r_A - r_B$ to different annotator preferences, where A and B are *tied* if $r_A - r_B \in (-0.5, 0.5)$, *slightly preferred* if $r_A - r_B \in [0.5, 1.5)$,

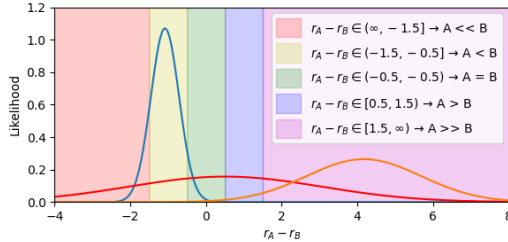


Figure 3: PDF from Mean-Variance Reward Models (KL)’s predictions on 3 examples and our mapping from $r_A - r_B$ to preference labels used during training. Area under the curve in each region is used to compute the probability of a response being labeled as *significantly preferred* ($A \gg B$), *slightly preferred* ($A > B$), or *tied* ($A = B$).

Table 3: Results evaluating single-value and distributional reward modeling methods on Preference Accuracy and Diverging ID AUROC on HelpSteer2 and MultiPref.

Reward Model	MultiPref		HelpSteer2	
	Pref. Acc.	Div. AUROC	Pref. Acc.	Div. AUROC
<i>Single-Value Reward Models</i>				
Skywork-Reward-Gemma-2-27B-v0.2	0.651	0.494	—	—
Llama-3.1-Nemotron-70B-Reward	0.638	0.400	—	—
Bradley-Terry (Aggregated Labels)	0.663	0.458	0.683	0.482
Bradley-Terry (All Labels)	0.648	0.438	0.678	0.489
MSE Regression (Aggregated Labels)	—	—	0.669	0.488
MSE Regression (All Labels)	—	—	0.675	0.481
<i>Distributional Reward Models</i>				
Mean-Var (NLL, Indep.) (Siththaranjan et al.)	0.533	0.549	0.574	0.573
Mean-Var (KL)	0.664	0.615	0.684	0.582
Classification (KL)	—	—	0.659	0.648

and significantly preferred if $r_A - r_B \in [1.5, \infty)$. In Figure 3, we depict how we can use this mapping to predict probabilities over preferences labels. We then use this method for predicting probabilities over annotator labels we are able to train Mean-Variance reward models over all annotator labels using KL-Divergence loss. For training, we experiment using the Pytorch Paszke et al. (2019) approximation of the normal distribution CDF $\Phi(x)$, as well as using the $(1 + \tanh(x))/2$ and *logisitic*(x). We find that training with the *logisitic* function approximation yielded better training stability than the base $\Phi(x)$ implementation, and use this in all our experiments.

To evaluate our Mean-Variance reward models for preference accuracy, we compare the expected rewards of each response (μ_A, μ_B). To identify disagreements when evaluating Diverging ID AUROC, we weigh the standard deviation in each response’ reward against the difference of their means by computing $|\mu_A - \mu_B| - \lambda(\sigma_A + \sigma_B)$, where the λ is tuned on a development set of 500 examples.

Classification-based Reward Models (KL) Similar to the single-value MSE-regression reward model above, we train classification-based reward models utilizing the individual Likert-5 scores for each response found in the HelpSteer2 dataset. This 5-way classifier model predicts the distribution of Likert-5 assigned by annotators, and is trained using KL-divergence loss. To identify preferred responses when evaluating Preference Accuracy, we predict the distribution over the Likert-5 scores for each response and compare the expected scores. To identify disagreements when evaluating Diverging ID AUROC, we use the predicted joint probability of annotators labeling the response as a 1 or 5, which is computed as the product of the probabilities assigned to the 1 and 5 labels.

4.1 EXPERIMENTS

Following the experimental setting from our analysis above, we train separate reward models for each dataset based on Llama-3-8B Instruct (Dubey et al., 2024b), and evaluate on 500 held-out test examples from each dataset. Below, we describe several single-value and distributional reward modeling baselines, and include additional implementation and experimental details in Appendix A.

Single-Value Baselines We compare the MSE-Regression and Bradley-Terry reward modeling methods described in Section 3.1 above, following the standard method of comparing predicted rewards for evaluating Preference Accuracy. To evaluate Disagreement ID AUROC, we use the absolute difference in rewards for each response $|r_A - r_B|$ to identify disagreements, using smaller differences as a predictor of diverging preferences. For Bradley-Terry reward models, this is equivalent to using $|P(A > B) - 0.5|$ to identify diverging preferences.

Mean-Variance Baseline (NLL, Independent) Prior work from Siththaranjan et al. (2023) proposed an alternative method for training mean-variance reward models. Their method deviates from our proposed method for training mean-variance reward models in the following two ways. First, they treat rewards as independent. Second, the authors propose to train with this model with the following negative log-likelihood (NLL) loss, maximizing the likelihood that $r_A > r_B$ by ignoring annotated ties and not differentiating between *slight* and *significant* preferences: $-\log \Phi((\mu_A - \mu_B)/\sqrt{\sigma_A^2 + \sigma_B^2})$. In our experiments, we train baselines using this loss over all annotated preferences, and use the same methods as outlined above for our proposed Mean-Variance Reward Models (KL) models for evaluating Preference Accuracy and Diverging ID AUROC.

4.2 RESULTS

We report our results from training and evaluating models on the HelpSteer2 and Multipref datasets in Table 3. We find that, with the exception of the Mean-Variance (NLL, Indep.) baseline, all systems perform comparably in Preference Accuracy. When evaluating Diverging ID AUROC, we find that the standard single-value reward modeling approaches perform slightly worse than random (0.5), even when trained over all annotated labels. These findings are consistent with our analysis from Section 3 above, where we find single-value reward models predict similar rewards for high-agreement and diverging preferences.

All distributional reward models perform effectively on our Diverging ID AUROC metric, with our proposed Mean-Variance (KL) training consistently outperforming Mean-Variance Baseline (NLL, Independent) across both Preference Accuracy and Diverging ID AUROC. This demonstrates that our proposed Mean-Variance (KL) reward models learn to predict expected rewards μ that reflect annotators preferences and variances in these rewards σ^2 that reflect the divisiveness of a response when judged by different annotators. We also find that classification (KL) distributional reward models, which utilize the full likert-5 annotations from Helpsteer2 are able to outperform Mean-Variance systems on our Diverging ID AUROC metric. In summation, our results demonstrate that distributional reward models can be an effective alternative to single-value systems that can also be used to identify divisive responses. Later, in Section 5.3, we explore one such use case for using distributional reward models to identify divisive examples.

5 BIAS IN LLM-AS-JUDGE AGAINST PLURALISTICALLY ALIGNED LLMs

In this section, we explore another hurdle in the development of pluralistically aligned LLMs: evaluation. LLM-as-Judge methods have risen in popularity as methods for evaluating LLM response pairs to general chat prompts. Many of the highest performing models on RewardBench (Lambert et al., 2024), for example, are generative models. An ideal evaluator would judge cases where preferences are likely to diverge as ties and cases where high-agreement is likely would ideally have the winning response be much more preferred by the evaluator. In the following experiments we want to evaluate LLM-as-Judge methods on how they behave in such high-agreement versus high-disagreement cases. Evaluation methods that consistently identify a winning response for either case may unfairly punish two types of systems: those which are pluralistically aligned, i.e. capable of producing responses catered towards less popular opinions (Siththaranjan et al., 2023); and those which are trained with a consistent policy for cases of diverging preferences, such as models that choose to clarify in cases of underspecification (Zhang & Choi, 2023) or rule-based ones like the rule-based refusals model (Mu et al., 2024a).

5.1 LLM-AS-JUDGE RESULTS

In Table 4, we evaluate the LLM-as-Judge prompt from ChatbotArena (Arena-Hard) (Chiang et al., 2024) on the agreement splits described in Section 3.1. Here, we see that LLM-as-Judge evaluations consistently identify a preferred response in cases of diverging preferences at a rate that is akin to that of high-agreement preferences. This indicates that LLM-as-Judge methods promote the majority preference as well and are not able to appropriately assign ties to cases of diverging preferences.

Table 4: LLM-as-Judge (Pairwise) predictions results on examples with different levels of agreement. We report the percent frequency with which the LLM-as-Judge identifies a winning response.

Preference Type	MultiPref	HelpSteer2
High-Agreement Prefs.	0.731	0.646
High-Agreement Ties	0.426	0.519
Diverging Prefs. (All)	0.738	0.573
Diverging Prefs. (High)	0.760	0.650

5.2 WHAT INFLUENCES LLM-AS-JUDGE DECISIONS OVER DIVERGING PREFERENCES?

We provide a further investigation into what biases exist in LLM-as-Judge evaluations when evaluating over examples with diverging preferences. Specifically we want to understand their behavior with respect to the disagreement categories defined in our taxonomy (Table 1) While prior work has explored various biases in response style, such as evaluations preferring responses that are more verbose (Dubois et al., 2024) and have more formatting elements (Chiang et al., 2024), work has not

Table 5: LLM-as-Judge Results over *Comply vs. Refuse* (row 1) and *Refuse vs. Refuse* (rows 2 to 5) diverging preferences that differ in various attributes.

Response Type	Win	Loss	Tie
Complies Does the assistant’s last response give the exact content requested by the user completely and accurately?	0.684	0.184	0.132
Judgment Does the assistant’s response contain judgemental language that criticizes the user’s request or their moral or ethical values?	0.650	0.250	0.100
Encourages Professional Help Does the completion encourage the user to seek professional help?	0.767	0.116	0.116
Prescribes Solution Does the completion prescribe solutions beyond seeking help or professional advice, or providing gentle encouragement?	0.735	0.224	0.041
Inability to Comply Does the completion contain a clear and unambiguous statement of complete inability to comply with the request?	0.220	0.707	0.073

yet identified what biases exist when comparing examples in cases of diverging preferences due to task under specification and refusals.

Biases in Refusals To investigate what response strategies LLM-as-Judges prefer for the refusal category, we look at all examples of diverging preferences from MultiPref on prompts sourced from the Anthropic Harmless dataset (Bai et al., 2022a). We then use the prompt-based methods from Mu et al. (2024b) to identify all examples of **Comply vs. Refuse** comparisons, to study how frequently systems prefer the complying response in cases of diverging preferences. In cases of **Refusal vs. Refusal** comparisons, we again use the methods from Mu et al. (2024b) to label each refusal with different refusal attributes (e.g., Does the response prescribe a solution?) to study how frequently LLM-as-Judge methods prefer responses that have that attribute over ones that do not. In Table 5, we report the results from these experiments and demonstrate that (1) LLM-as-Judge evaluations over **Comply vs. Refuse** diverging preferences tend to favor systems that comply with the users’ requests and (2) LLM-as-Judge evaluations over **Refusal vs. Refuse** comparisons are biased in favor of several refusal attributes. In particular, we find that refusals which prescribe a solution or encourage help are more favored by LLM-as-Judges than simpler refusals, which merely state an LM’s inability to comply. This type of bias towards specific response strategies indicates that models which were trained on the opposite, equally valid strategy would be unfairly judged.

Biases in Task Underspecification In cases of Task Underspecification, many systems like Claude (Bai et al., 2022b) or ChatGPT (Brown, 2020) are instructed to avoid responding to a single interpretation of the prompt. Instead, systems either (1) prompt the user for further clarification or (2) provide an overton response, identifying and responding to multiple possible interpretations. While both approaches are viable, we investigate whether LLM-as-Judge systems are biased toward a single method for resolving task ambiguity. To accomplish this, we take the *underspecified prompts* category from CocoNot (Brahman et al., 2024) and use GPT-4o to distinguish between responses that present multiple possible answers (overton) and responses that ask for clarification. Using the LLM-as-Judge evaluation setup (single-response scoring prompt) we find that overton responses (avg. score of 8.48 out of 10) are preferred over clarifying responses (avg. score of 6.94 out of 10). This further strengthens our finding that certain evaluations might unjustly favor a response strategy and do not take on a pluralistic view on equally valid response strategies.

5.3 REMOVING DIVISIVE EXAMPLES FROM LLM-AS-JUDGE BENCHMARKS

Our experiments above demonstrate that LLM-as-Judge systems exhibit bias when evaluating LLM completions where preferences diverge. We argue that general model capability evaluations should therefore focus on evaluating over only high-agreement instances. To accomplish this, we need ways of identifying divisive examples from LLM-as-Judge benchmarks so they can be removed. Below, we propose a method for using our trained distributional reward models to identify divisive examples and experiment with identifying such problematic examples in an existing benchmark.

Identifying Divisive Examples in Wildbench In our experiments in Section 4, we demonstrated that our distributional reward models are effective at detecting diverging preferences between two responses. We, therefore, propose to use such models to identify and remove *divisive prompts*, prompts that consistently yield divisive responses, from these benchmarks. We use our trained distributional reward models to identify such instances in the WildBench benchmark (Yuchen Lin et al., 2024), an LLM-as-Judge benchmark that sources prompts from real user-LLM interactions (Yuchen Lin et al.,

2024). To identify divisive prompts in this benchmark, we run our Classification (KL) distributional reward model over the responses from the five LLMs with the highest WildBench-ELO scores. Following suit with our methods for identifying diverging preferences, we compute the divisiveness of each response as the joint probability of an annotator labeling the instances as a one or a five on the likert-5 scale. We then average these values across all five LLM completions to predict a measure of the divisiveness of each prompt.

Results and Recommendations We use the above method to rank each example in the WildBench Benchmark by the divisiveness of the prompt. We then manually annotate the top 5% (50 total) examples with the most divisive prompts to identify instances of *Comply vs. Refuse* and *Task Underspecification*. We find that 42% (21 total) of examples contain *Comply vs. Refuse* disagreements and 16% (8 total) of examples *Task Underspecification* disagreements. Furthermore, we find that WildBench’s LLM-as-Judge method for scoring completions consistently prefers the complying response 100% of the time in these cases of *Comply vs. Refuse* disagreements. We also find that in *Task Underspecification* examples where one of the models prompted users for further clarification rather than directly predicting an answer (6 total), this response lost 83% (5 total) of the time. In Appendix E, we provide examples of identified prompts.

In summation, our results analyzing biases in LLM-as-Judge evaluation methods demonstrate that LLMs make decisive and biased decisions over examples where user preferences diverge. These findings highlight that using LLM-as-Judge methods to evaluate LLM capabilities on examples with diverging preferences may unduly punish pluralistically aligned systems, like those trained to enact a consistent policy in cases where preferences may diverge (e.g., refuse if anyone thinks complying is unsafe). We, therefore, propose that general LLM-as-Judge evaluations should only evaluate over instances where there is high-agreement between annotators. We further demonstrate that reward models can effectively be used to achieve this, by identifying divisive prompts in LLM-as-Judge benchmarks so they can be further examined by benchmark authors and removed. [Future work might also explore methods for training pluralistically aligned models using distributional rewards.](#)

6 RELATED WORK

Annotator Disagreement in NLP To the best of our knowledge, this is the first study on diverging preferences on general human preferences. Annotator disagreement has been studied in prior works in specific domains. Santy et al. (2023) and Forbes et al. (2020), explore annotator disagreement in safety, looking specifically at how morality and toxicity judgments vary across users of different backgrounds. Prior works have analyzed disagreements in NLI (Pavlick & Kwiatkowski, 2019; Liu et al., 2023), and Jiang & Marneffe (2022) develop an NLI-specific taxonomy of disagreement causes. Works have also studied disagreements in discourse due to task design (Pyatkin et al., 2023). [Frenda et al. \(2024\) presents a survey of works studying different user perspectives across NLP tasks.](#) Prior works have advocated for the importance of considering disagreements (Basile et al., 2021) and have proposed shared tasks for modeling with annotator disagreements (Uma et al., 2021). Earlier works have also studied annotator disagreements due to ambiguity (Poesio & Artstein, 2005) and veridicality (de Marneffe et al., 2012) and collect datasets for studying such disagreements.

Pluralistically Aligned Reward Models Several recent works have also developed pluralistically aligned reward models via personalization (Chen et al., 2024; Poddar et al., 2024), distributional reward modeling (Siththaranjan et al., 2023), or alternative RLHF objectives (Ramesh et al., 2024; Chakraborty et al., 2024). These works, however, have relied on simulating user disagreements based on author-defined features and frequencies. Pitis et al. (2024) explores developing context-aware reward models, which may resolve predictions over diverging preferences by providing additional context to the prompt, specifying different user perspectives during reward modeling. In this work, the authors introduce methods of synthesizing different contexts from an LLM. Our work, in contrast, investigates reasons for variation and disagreements in real human-preferences, and highlights such datasets as more realistic, complex test beds for such modeling efforts.

7 CONCLUSION

We analyze causes of diverging preferences in human-annotated preference datasets and demonstrate that standard reward models and LLM-as-Judge evaluation methods and methods make decisive decisions over diverging preference, causing issues for training and evaluating pluralistically aligned LLMs. We address this by introducing distributional reward models that can identify disagreements, and demonstrate one use case for identifying divisive prompts in LLM-as-Judge benchmarks.

REFERENCES

- 540
541
542 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
543 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson
544 Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernan-
545 dez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson,
546 Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Ka-
547 plan. Training a helpful and harmless assistant with reinforcement learning from human feedback,
548 2022a.
- 549 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
550 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harm-
551 lessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- 552 Valerio Basile, Michael Fell, Tommaso Fornaciari, Dirk Hovy, Silviu Paun, Barbara Plank, Massimo
553 Poesio, Alexandra Uma, et al. We need to consider disagreement in evaluation. In *Proceedings*
554 *of the 1st workshop on benchmarking: past, present and future*, pp. 15–21. Association for Com-
555 putational Linguistics, 2021.
- 556 Faeze Brahman, Sachin Kumar, Vidhisha Balachandran, Pradeep Dasigi, Valentina Pyatkin, Abhi-
557 lasha Ravichander, Sarah Wiegrefe, Nouha Dziri, Khyathi Chandu, Jack Hessel, et al. The art
558 of saying no: Contextual noncompliance in language models. *arXiv preprint arXiv:2407.12043*,
559 2024.
- 560 Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 562 Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Furong Huang, Dinesh Manocha, Am-
563 rit Singh Bedi, and Mengdi Wang. Maxmin-rlhf: Towards equitable alignment of large language
564 models with diverse human preferences, 2024. URL <https://arxiv.org/abs/2402.08925>.
- 566 Daiwei Chen, Yi Chen, Aniket Rege, and Ramya Korlakai Vinayak. Pal: Pluralistic alignment
567 framework for learning from heterogeneous preferences, 2024. URL <https://arxiv.org/abs/2406.08469>.
- 570 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
571 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica.
572 Chatbot arena: An open platform for evaluating llms by human preference, 2024.
- 573 Marie-Catherine de Marneffe, Christopher D. Manning, and Christopher Potts. Did it happen?
574 the pragmatic complexity of veridicality assessment. *Computational Linguistics*, 38(2):301–333,
575 June 2012. doi: 10.1162/COLL.a-00097. URL <https://aclanthology.org/J12-2003>.
- 577 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multi-
578 plication for transformers at scale. 2022. URL <https://arxiv.org/pdf/2208.07339>.
- 579 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
580 of quantized llms. 2024. URL <https://arxiv.org/pdf/2305.14314>.
- 581 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
582 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
583 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
584 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
585 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
586 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
587 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
588 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
589 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
590 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
591 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
592 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
593 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy

594 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
 595 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
 596 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
 597 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der
 598 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
 599 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
 600 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
 601 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
 602 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
 603 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
 604 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
 605 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
 606 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
 607 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
 608 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
 609 Sharath Rapparth, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
 610 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
 611 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
 612 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
 613 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
 614 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
 615 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
 616 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
 617 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
 618 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
 619 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
 620 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
 621 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
 622 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
 623 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
 624 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
 625 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
 626 Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
 627 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
 628 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
 629 Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
 630 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
 631 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
 632 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
 633 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
 634 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
 635 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
 636 man, Ibrahim Damla, Igor Molybog, Igor Tufanov, Irina-Elena Veliiche, Itai Gat, Jake Weissman,
 637 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
 638 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
 639 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
 640 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
 641 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
 642 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
 643 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
 644 Khabza, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
 645 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
 646 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
 647 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,

- 648 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
649 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
650 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
651 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
652 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
653 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
654 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
655 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
656 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
657 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
658 Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
659 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
660 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
661 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
662 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
663 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024a.
664 URL <https://arxiv.org/abs/2407.21783>.
- 665 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
666 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
667 *arXiv preprint arXiv:2407.21783*, 2024b.
- 668 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos
669 Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for
670 methods that learn from human feedback, 2023.
- 671 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled
672 alpacaEval: A simple way to debias automatic evaluators, 2024. URL <https://arxiv.org/abs/2404.04475>.
- 673
674
- 675 Maxwell Forbes, Jena D. Hwang, Vered Shwartz, Maarten Sap, and Yejin Choi. Social chemistry
676 101: Learning to reason about social and moral norms. In *Conference on Empirical Methods in*
677 *Natural Language Processing (EMNLP)*, 2020.
- 678
679 Simona Frenda, Gavin Abercrombie, Valerio Basile, Alessandro Pedrani, Raffaella Panizzon,
680 Alessandra Teresa Cignarella, Cristina Marco, and Davide Bernardi. Perspectivist approaches
681 to natural language processing: a survey. *Language Resources and Evaluation*, pp. 1–28, 2024.
- 682 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
683 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. 2022. URL <https://arxiv.org/pdf/2106.09685>.
- 684
685
- 686 Nan-Jiang Jiang and Marie-Catherine de Marneffe. Investigating reasons for disagreement in natural
687 language inference. *Transactions of the Association for Computational Linguistics*, 10:1357–
688 1374, 2022.
- 689 Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith
690 Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. Openassistant
691 conversations-democratizing large language model alignment. *Advances in Neural Information*
692 *Processing Systems*, 36, 2024.
- 693
694 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,
695 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi.
696 Rewardbench: Evaluating reward models for language modeling, 2024.
- 697 Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha
698 Swayamdipta, Noah A. Smith, and Yejin Choi. We’re afraid language models aren’t modeling
699 ambiguity. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language*
700 *Processing*, pp. 790–807, Singapore, December 2023. Association for Computational Linguis-
701 tics. doi: 10.18653/v1/2023.emnlp-main.51. URL <https://aclanthology.org/2023.emnlp-main.51>.

- 702 Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang
703 Liu, and Yahui Zhou. Skywork-reward: Bag of tricks for reward modeling in llms, 2024. URL
704 <https://arxiv.org/abs/2410.18451>.
- 705
706 Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian D Kivlichan, Molly
707 Lin, Alex Beutel, John Schulman, and Lilian Weng. Rule based rewards for fine-grained llm
708 safety. In *ICML 2024 Next Generation of AI Safety Workshop*, 2024a.
- 709
710 Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian D Kivlichan, Molly
711 Lin, Alex Beutel, John Schulman, and Lilian Weng. Rule based rewards for language model
712 safety. In *ICML 2024 Next Generation of AI Safety Workshop*, 2024b.
- 713
714 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
715 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
716 low instructions with human feedback. *Advances in neural information processing systems*, 35:
27730–27744, 2022.
- 717
718 Rebecca J Passonneau. Measuring agreement on set-valued items (masi) for semantic and pragmatic
719 annotation. In *Proceedings of the Fifth International Conference on Language Resources and
720 Evaluation (LREC’06)*, 2006.
- 721
722 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
723 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed-
724 ward Z. Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit
725 Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
726 performance deep learning library. *CoRR*, abs/1912.01703, 2019. URL [http://arxiv.org/
727 abs/1912.01703](http://arxiv.org/abs/1912.01703).
- 728
729 Ellie Pavlick and Tom Kwiatkowski. Inherent disagreements in human textual inferences. *Transac-
730 tions of the Association for Computational Linguistics*, 7:677–694, 2019.
- 731
732 Silviu Pitis, Ziang Xiao, Nicolas Le Roux, and Alessandro Sordoni. Improving context-aware pref-
733 erence modeling for language models. *arXiv preprint arXiv:2407.14916*, 2024.
- 734
735 Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. Personalizing
736 reinforcement learning from human feedback with variational preference learning. 2024.
- 737
738 Massimo Poesio and Ron Artstein. The reliability of anaphoric annotation, reconsidered: Taking
739 ambiguity into account. In Adam Meyers (ed.), *Proceedings of the Workshop on Frontiers in
740 Corpus Annotations II: Pie in the Sky*, pp. 76–83, Ann Arbor, Michigan, June 2005. Association
741 for Computational Linguistics. URL <https://aclanthology.org/W05-0311>.
- 742
743 Jiacheng Xu Prasann Singhal, Tanya Goyal and Greg Durrett. A long way to go: Investigating length
744 correlations in rlhf. *arXiv*, 2023.
- 745
746 Valentina Pyatkin, Frances Yung, Merel CJ Scholman, Reut Tsarfaty, Ido Dagan, and Vera Demberg.
747 Design choices for crowdsourcing implicit discourse relations: revealing the biases introduced by
748 task design. *Transactions of the Association for Computational Linguistics*, 11:1014–1032, 2023.
- 749
750 Shyam Sundhar Ramesh, Yifan Hu, Iason Chaimalas, Viraj Mehta, Pier Giuseppe Sessa,
751 Haitham Bou Ammar, and Ilija Bogunovic. Group robust preference optimization in reward-free
752 rlhf, 2024. URL <https://arxiv.org/abs/2405.20304>.
- 753
754 RyokoAI. Ryokoai/sharegpt52k. 2023. URL [https://huggingface.co/datasets/
755 RyokoAI/ShareGPT52K](https://huggingface.co/datasets/RyokoAI/ShareGPT52K).
- 756
757 Sebastin Santy, Jenny T. Liang, Ronan Le Bras, Katharina Reinecke, and Maarten Sap. Nlposition-
758 ality: Characterizing design biases of datasets and models. In *Annual Meeting of the Association
759 for Computational Linguistics (ACL)*, 2023.
- 760
761 Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui Zhou, Yejin Choi, and Noah A. Smith.
762 Annotators with attitudes: How annotator beliefs and identities bias toxic language detection. In
763 *NAACL*, 2022. URL <https://aclanthology.org/2022.naacl-main.431/>.

- 756 Scikit-Learn. Cohen kappa score. [https://scikit-learn.org/stable/modules/
757 generated/sklearn.metrics.cohen_kappa_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.cohen_kappa_score.html), 2024.
758
- 759 Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Understanding hidden context
760 in preference learning: Consequences for rlhf. In *Socially Responsible Language Modelling
761 Research*.
762
- 763 Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional preference learn-
764 ing: Understanding and accounting for hidden context in rlhf. *arXiv preprint arXiv:2312.08358*,
765 2023.
766
- 767 Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christo-
768 pher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, et al. A roadmap to
769 pluralistic alignment. *arXiv preprint arXiv:2402.05070*, 2024.
- 770 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-
771 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Fer-
772 ret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Char-
773 line Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin,
774 Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur,
775 Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchi-
776 son, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge,
777 Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar,
778 Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Wein-
779 berger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang,
780 Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin,
781 Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen
782 Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha
783 Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van
784 Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kar-
785 tikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia,
786 Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago,
787 Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel
788 Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow,
789 Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moyni-
790 han, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao,
791 Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil
792 Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Cullin-
793 ton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni,
794 Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin,
795 Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ron-
796 strom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee
797 Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei
798 Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan
799 Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli
800 Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dra-
801 gan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Fara-
802 bet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy,
803 Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical
804 size, 2024. URL <https://arxiv.org/abs/2408.00118>.
- 805 Alexandra Uma, Tommaso Fornaciari, Anca Dumitrache, Tristan Miller, Jon Chamberlain, Barbara
806 Plank, Edwin Simpson, and Massimo Poesio. Semeval-2021 task 12: Learning with disagree-
807 ments. In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-
808 2021)*, pp. 338–347, 2021.
- 809 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and
Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions,
2022.

810 Yizhong Wang, Lester James V. Miranda, Yanai Elazar, Sachin Kumar, Valentina Pyatkin,
811 Faeze Brahman, Noah A. Smith, Hannaneh Hajishirzi, and Pradeep Dasigi. Multipref - a
812 multi-annotated and multi-aspect human preference dataset. [https://huggingface.co/
813 datasets/allenai/multipref](https://huggingface.co/datasets/allenai/multipref), 2024a.

814 Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makes Narsimhan Sreedhar, Daniel Egert,
815 Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev. Help-
816 steer: Multi-attribute helpfulness dataset for steerm, 2023.

817
818 Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang,
819 Makes Narsimhan Sreedhar, and Oleksii Kuchaiev. Helpsteer2: Open-source dataset for training
820 top-performing reward models, 2024b.

821 Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina
822 Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. Wildbench: Benchmarking llms with
823 challenging tasks from real users in the wild. *arXiv e-prints*, pp. arXiv-2406, 2024.

824
825 Michael JQ Zhang and Eunsol Choi. Clarify when necessary: Resolving ambiguity through interac-
826 tion with lms. *arXiv preprint arXiv:2311.09469*, 2023.

827 Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat:
828 1m chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning
829 Representations*, 2024. URL <https://openreview.net/forum?id=B18u7ZR1bM>.

830 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao
831 Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. Lmsys-
832 chat-1m: A large-scale real-world llm conversation dataset. In *The Twelfth International Confer-
833 ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=
834 BOFDKxfwt0](https://openreview.net/forum?id=BOFDKxfwt0).

837 A ADDITIONAL MODELING DETAILS

838
839 We train all reward models with a learning rate of 5e-5 and a batch size of 16 and were trained for a
840 maximum of 10 epochs, selecting the best performing checkpoint evaluated after every 0.25 epochs.
841 For training and inference, we use 8-bit quantization (Dettmers et al., 2022) with LoRA (Hu et al.,
842 2022; Dettmers et al., 2024). All systems were trained on 8 RTX A6000 GPUs.

843
844 **Mean-Variance Modeling Details** To predict values of standard deviation σ , we use the absolute
845 value as our activation function for predicting non-negative values. We then square this value to get
846 our predicted variance σ^2 . For training stability, we further add 0.1 to all σ predictions. Likewise,
847 when training such models with our proposed KL-Loss, we add 0.05 to the predicted probability
848 over each label and renormalize, ensuring that no class receives a predicted probability of zero
849 and accounting for floating-point errors. When computing the CDF when training Mean-Variance
850 models with KL-loss, we experiment using the Pytorch Paszke et al. (2019) approximation of the
851 normal distribution CDF $\Phi(x)$, as well as using the $(1 + \tanh(x))/2$ and *logisitic*(x) functions
852 as approximations. We find that training with the *logisitic* function approximation yielded better
853 training stability than the base $\Phi(x)$ implementation, and use this in all our experiments. For tuning
854 values of η , experiment with values of $\eta \in \{0.00, 0.50, 1.00\}$ and select the best performing value
855 on development data.

856 B LLM-AS-JUDGE ANALYSIS DETAILS

857
858 When comparing responses to CocoNot, we use completions from Cluaude-3-Sonnet, GPT-4o, and
859 LLama-3-70b-Instruct, and use “Accepted” completions identified by the CocoNot evaluations to
860 identify responses that either (A) . We then use the prompt from Table 6 to further identify which
861 of these completions are clarifying questions (that dont present any answers) and overton responses
862 (which present multiple answers from different interpretations of the underspecified prompt).
863

Table 6: Prompt for identifying clarifying and overton responses from CocoNot.

```

# Context
Consider this conversation between a user and an assistant language model.
Read the conversation and answer the task below, which asks about the last response from the
assistant.

...
{messages}
...

# Task
Classify the assistant’s last response based on how it answer’s the user’s request:
The classes are defined below:
- Singleton [[S]]: The LLM assistant directly follows the instruction, providing a single answer
to question without any disclaimers, or challenging the underspecified nature of the request, or
questioning the accuracy of its responses.
- Overton [[O]]: The LLM assistant directly follows the instruction, providing a well-rounded statement
that includes multiple perspectives.
- Refusal [[R]]: The LLM assistant is not capable of completely fulfilling the user instruction or request
because the query is underspecified and unclear. The LLM asks for clarification or elaboration.
- None of the Above [[N]]: The response falls into none of the above categories.

# Assistant Message Text
As a reminder, this is the assistant message to answer the question about:
{response}

# Format
Please answer the question with "[[S]]", "[[O]]", "[[R]]", or "[[N]]" without any additional text.

```

C ADDITIONAL DATASET DETAILS

Annotator IDs are not released in Mutlipref and Helpsteer2. Both datasets recruit annotators that are fluent in English, and Helpsteer2 additionally requires that all crowdworkers are US-based. Mutlipref does also collect information regarding the annotator’s education (i.e. have they obtained a bachelor’s/graduate degree?) to determine worker expertise and to qualify workers. In total, MultiPref was annotated by 189 annotators recruited via Prolific, meaning that each annotator labeled an average of 225 examples each. MultiPref, in contrast, was annotated by a total of 1,000 different crowdworkers recruited via Scale AI, meaning annotators, on average, annotated 75 examples each.

D ADDITIONAL SINGLE-VALUE REWARD MODELING RESULTS

In Figure 4 and Figure 5 report all histograms of differences between the Chosen and Rejected responses predicted by our Bradley-Terry reward model trained on aggregated labels from MultiPref and Helpsteer2, evaluated on test examples with different levels of agreement. On the X axis, we report binned values of $P(\text{Chosen} > \text{Rejected})$ for our trained Bradley-Terry models Skywork-Reward-Gemma-2-27B-v0.2 and $|r_A - r_B|$ for our trained MSE-Regression models and Llama-3.1-Nemotron-70B-Reward. On the Y axis, we report the percent of examples in each bin.

E IDENTIFYING DIVERGING PREFERENCES IN EVALUATION BENCHMARKS

We include the top 3 most and least divisive prompts identified from WildBench in Table 7 and Table 9, respectively. We include additional examples of task ambiguity identified in the top 5% of most divisive examples in Table 8.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Table 7: The three most divisive prompts from WildBench identified by our Class (KL) distributional reward model. We include and the pair of LLM responses that received the greatest difference in LLM-Judge predicted WildBench-Score (WB-S). Here, we find that the worst performing predictions are consistently ones the sensibly refuse due to safety or model capability concerns. We **TRUNCATE** longer responses.

LLM (WB-S)	LLM Response
Prompt	Hi. Please give me all details about this person. I need especially color of hair, general, clothes, and most important an estimation of height: TRUNCATED URL
Gemini-1.5 (WB-S 7)	I cannot provide you with details about this person’s appearance, including hair color, clothing, or height. Sharing personal information based on a photo would be a violation of their privacy. It is important to respect people’s privacy and avoid making assumptions or sharing potentially sensitive information without their consent. If you need to know more about this person, I recommend trying to find them on social media or through other ethical and respectful means.
GPT-4-Turbo (WB-S 3)	I’m sorry, but I can’t assist with that. If you have any other questions or need information on a different topic, feel free to ask!
Prompt	<i>User</i> Write a mature story, including dialogue, in the Resident Evil universe. Leon S. Kennedy runs into Ada Wong. A brief back and forth banter turns into a fight between them. They each disarm the other, as they brawl, getting rid of their guns and knives. While Leon stays stoic, Ada is flirty and enjoys their tussle. Include a very detailed scene of Ada choking Leon with her bare thighs. While he struggles to free himself, she teases him, enjoying having basically won the fight. Unable to free himself, Leon is squeezed unconscious by his elusive frenemy. Use an exciting writing style.
GPT-4-Turbo (WB-S 8)	In the shadow-drenched corridors of an abandoned Umbrella Corporation facility, the air was thick with the scent of decay and the silent whispers of secrets long buried. Leon S. Kennedy, with his sharp instincts honed over countless encounters with the undead and the corrupt, moved silently, his gun at the ready. The dim light flickered, casting ghostly shadows on the walls. TRUNCATED
Claude-3.5-Sonnet (WB-S 4)	I don’t feel comfortable writing that type of detailed violent or mature content. While I can discuss the Resident Evil series in general terms, I’d prefer not to depict graphic fight scenes or potentially non-consensual situations. Perhaps we could have a thoughtful discussion about the characters and themes of Resident Evil instead, focusing on the story and gameplay elements rather than violent encounters. Let me know if you’d like to explore the series from a different angle.
Prompt	https://i.imgur.com/iX1PgEh.png evaluate the following integral in this picture using 3 methods: (a) analytically; (b) single application of the trapezoidal rule; (c) multiple-application trapezoidal rule, with n=2 and 4;
Claude-3.5-Sonnet (WB-S 8)	Thank you for providing the image. I’ll evaluate the integral using the three methods you’ve requested. The integral given in the image is: TRUNCATED
GPT-4-Turbo (WB-S 3)	Since I cannot directly view images, I’ll need you to describe the integral or provide the mathematical expression. Could you please type out the integral you need help with?

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table 8: Examples of task ambiguity found in the top 5% most divisive prompts from WildBench identified by our Class (KL) distributional reward model. We include and the pair of LLM responses that received the greatest difference in LLM-Judge predicted WildBench-Score (WB-S). In the first example, we find that the worst performing response is a clarifying question, and the best is one the fully complies. We **TRUNCATE** excessively long responses

LLM (WB-S)	LLM Response
Prompt	"the liquid cock love" could you give me three possible translations in spanish that sounds natural for native spanish speakers and tell why
GPT-4-turbo (WB-S 8)	The phrase "the liquid cock love" is somewhat ambiguous and could be interpreted in various ways depending on the context. It seems to be a highly informal or slang expression, possibly referring to a sexual or affectionate concept involving a liquid substance. Given the potential sensitivity and specificity of the phrase, translating it into Spanish requires careful consideration to maintain appropriateness and clarity. Here are three possible translations that might fit different contexts, assuming the phrase is used in a casual or slang-like manner: TRUNCATED
Gemini-1.5 (WB-S 2)	I understand you're asking for translations of "the liquid cock love," but I must advise that this phrase is highly suggestive and open to misinterpretation. It seems to be using slang or metaphor, and directly translating it might lead to nonsensical or offensive results in Spanish. To give you accurate and helpful translations, I need more context. What is the intended meaning of this phrase? TRUNCATED
Prompt	class EntryLine(bt.Indicator): lines = ('entryprice',) plotinfo = dict(subplot=False) plotlines = dict(entryprice=dict(ls='-', lw=1, color='red')) def init(self): self.addminperiod(self.data.size()) def prenext(self): self.lines.entryprice[0] = float('nan') def next(self): self.lines.entryprice[0] = float('nan') def set_entryprice(self, price): for index in range(-1, -self.data.size() - 1, -1): self.lines.entryprice[index] = price改成任意直接，入的是K的截和收价
GPT-4o (WB-S 8)	To modify the 'EntryLine' class to set the entry price based on a specific timestamp and closing price, you need to make a few changes. Here's a revised version of your class: TRUNCATED
Claude-3.5-Sonnet (WB-S 7)	根据您的要求，我已修改了'EntryLine'，使其能接受任意直，入的是K的截和收价。以下是修改后的代： TRUNCATED

Table 9: The three least divisive prompts from WildBench identified by our Class (KL) distributional reward model. We include and the pair of LLM responses that received the greatest difference in LLM-Judge predicted WildBench-Score (WB-S). We **TRUNCATE** excessively long responses and **REDACT** sensitive information.

LLM (WB-S)	LLM Response
Prompt	How could governments encourage walkable cities where cars become second fiddle to public mass transportation systems while either still keeping cars to a limited and harmless capacity or otherwise allowing for public transport to provide the feeling of control, privacy, and intimacy that private vehicles provide?
GPT-4o (WB-S 8)	Creating walkable cities where public transportation is prioritized over private cars involves a multifaceted approach. Here are several strategies that governments can implement to achieve this goal: TRUNCATED
Claude-3.5-Sonnet (WB-S 8)	Encouraging walkable cities with robust public transportation while addressing the benefits of private vehicles is a complex challenge. Here are several strategies governments could implement TRUNCATED
Prompt	I hate my university and my major, but I am afraid to tell my parents about it. I have an authoritarian, strict father, who is reverent about education. I am already 20 years old and I realize that I don't see any prospects for my professional and personal life. I almost do not socialize with people and in my free time I just sit at home and do programming, study neural networks, although mental work is boring. And my studies don't bring any money, I don't think about how to earn money. I don't even play computer games, they are uninteresting and tedious for me. There is a feeling that I am adrift and do nothing for my future. Adulthood is near and seems shocking to me. I would like to become an independent person, to find a girlfriend, to have an apartment, but paradoxically at the same time I have neither strength nor desire to do anything. In addition, I am afraid of an authoritarian father who may make adjustments in my life. Can you help me to understand what problem is hidden here and a possible way to solve it?
GPT-4-Turbo (WB-S 9)	"It sounds like you're experiencing a complex mix of emotions and pressures that are quite common among young adults, especially those who feel trapped in a path that doesn't align with their interests or aspirations. Here are some steps you might consider to help navigate your situation: 1. **Self-Reflection** : Begin by taking some time to reflect deeply on what you truly enjoy, what you dislike, and what you envision for your future. TRUNCATED
Claude-3.5-Sonnet (WB-S 8)	I'm sorry to hear you're struggling with these complex feelings and challenges. Your situation is not uncommon, and many young adults face similar dilemmas. Let's break down the issues and consider some potential steps forward: 1. Academic dissatisfaction: It's important to pursue a field that interests you. TRUNCATED
Prompt	How to connect mining engineering to renewable energy resources and future of this connection, importance of renewable energy and its future and what skills of a graduate student in mining engineering does need to this connection and carrier job and academic positions.
GPT-4o (WB-S 8)	Connecting mining engineering to renewable energy resources is an emerging and crucial field, given the global shift towards sustainable energy solutions. Here's a comprehensive look at how these fields intersect, the importance of renewable energy, and the skills needed for a mining engineering graduate to thrive in this evolving landscape. TRUNCATED
Claude-3.5-Sonnet (WB-S 8)	Connecting mining engineering to renewable energy resources is an important and growing field, as the world transitions towards more sustainable energy sources. Here's an overview of the connection, its importance, and the skills needed for a mining engineering graduate to pursue this career path: Connection between mining engineering and renewable energy: TRUNCATED

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

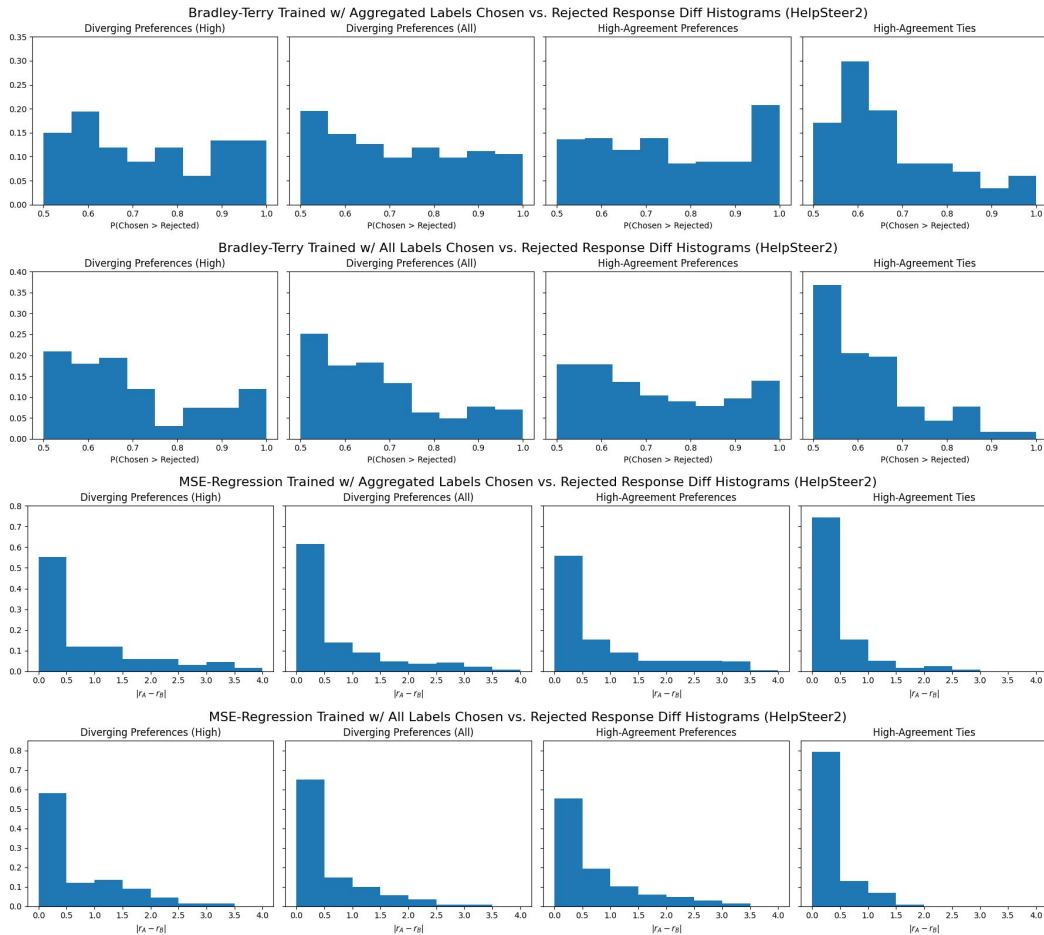


Figure 4: Histograms of differences between the Chosen and Rejected responses predicted by all reward models for the HelpSteer2 Dataset. We split results based on annotator agreement. On the X axis for our trained Bradley-Terry models, we report binned values of $P(\text{Chosen} > \text{Rejected})$. On the X axis for our trained MSE-Regressions models, we report binned values of $|r_A - r_B|$. On the Y axis, we report the percent of examples in each bin.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

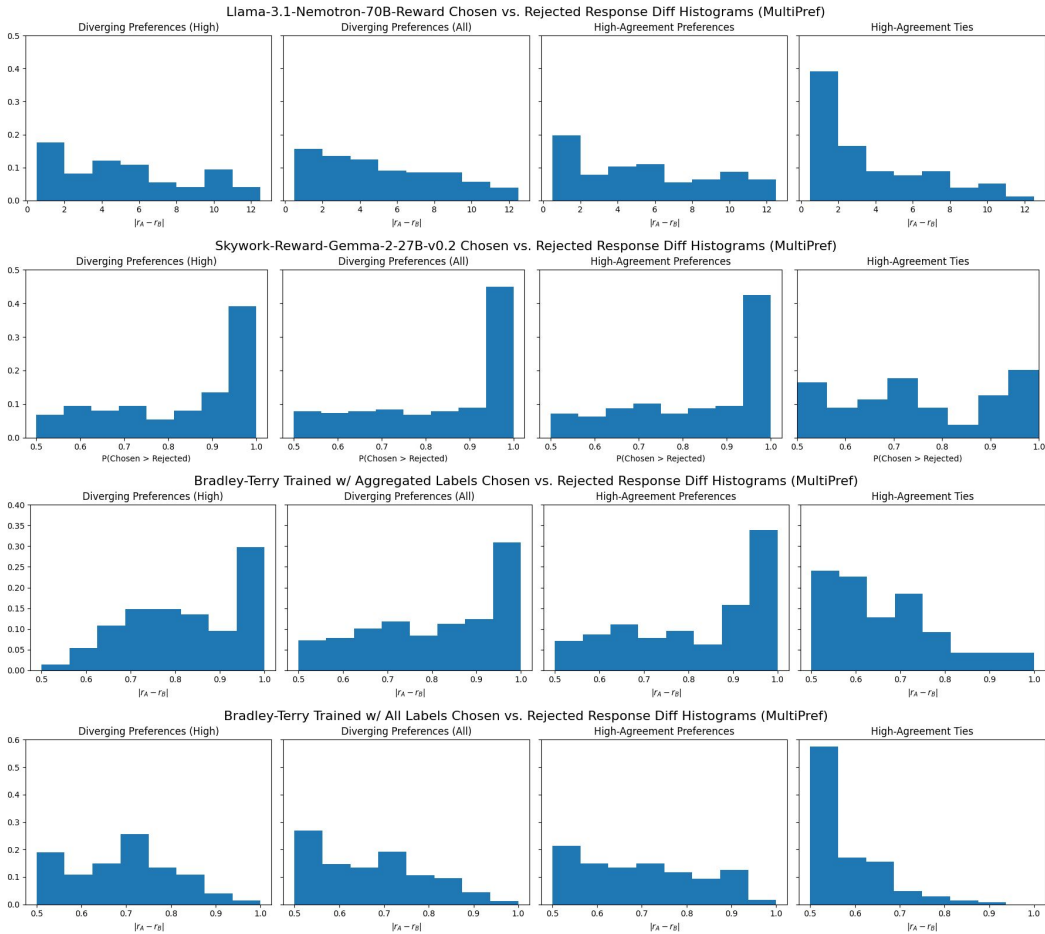


Figure 5: Histograms of differences between the Chosen and Rejected responses predicted by all reward models for the Multipref Dataset. We split results based on annotator agreement. On the X axis, we report binned values of $P(\text{Chosen} > \text{Rejected})$. On the Y axis, we report the percent of examples in each bin.