

WildFeedback: Aligning LLMs With In-situ User Interactions And Feedback

Anonymous ACL submission

Abstract

As large language models (LLMs) continue to advance, aligning these models with human preferences has emerged as a critical challenge. Traditional alignment methods, relying on human or LLM annotated datasets, are limited by their resource-intensive nature, inherent subjectivity, misalignment with real-world user preferences, and the risk of feedback loops that amplify model biases. To overcome these limitations, we introduce WILDFEEDBACK, a novel framework that leverages in-situ user feedback during conversations with LLMs to create preference datasets automatically. Given a corpus of multi-turn user-LLM conversation, WILDFEEDBACK identifies and classifies user feedback to LLM responses between conversation turns. The user feedback is then used to create examples of preferred and dispreferred responses according to users' preference. Our experiments demonstrate that LLMs fine-tuned on WILDFEEDBACK dataset exhibit significantly improved alignment with user preferences, as evidenced by both traditional benchmarks and our proposed checklist-guided evaluation. By incorporating in-situ feedback from actual users, WILDFEEDBACK addresses the scalability, subjectivity, and bias challenges that plague existing approaches, marking a significant step toward developing LLMs that are more responsive to the diverse and evolving needs of their users.

1 Introduction

Large language models (LLMs) have become a cornerstone of modern natural language processing (NLP) applications, powering a wide range of tasks from conversational agents to content generation. Despite their strengths, aligning LLMs with human preferences remains a challenge (Bai et al., 2022a; Ouyang et al., 2022; OpenAI et al., 2024; Dubey et al., 2024). Traditional alignment methods involve instruction tuning and preference training on curated human or LLM-annotated datasets (Bai

et al., 2022a; Ouyang et al., 2022; Cui et al., 2024). However, these approaches face critical limitations: human annotation is resource-intensive and often subjective, while LLM-generated synthetic data risks reinforcing biases instead of capturing diverse human preferences (Gautam and Srinath, 2024; Wyllie et al., 2024; Chen et al., 2024; Poddar et al., 2024).

In response, recent work explores in-situ user feedback (e.g., upvotes, downvotes, engagement) for LLM training (Shi et al., 2022; Lin et al., 2024b; Don-Yehiya et al., 2024). This approach harnesses authentic user feedback during interactions with LLMs, offering a more dynamic and accurate reflection of user preferences. Rather than relying on static, costly, and misaligned pre-collected data, this method adapts to evolving user needs. However, existing works are limited in scope, either requiring explicit, structured feedback from users or fine-tuning models directly on responses that trigger explicit user feedback.

In this paper, we introduce WILDFEEDBACK, a novel framework designed to align LLMs with in-situ user interactions and feedback. WILDFEEDBACK addresses the limitations of existing approaches by constructing preference datasets from real user-LLM conversations, specifically focusing on user feedback that naturally occurs during these interactions. The overview of the framework is shown in Figure 1. Our framework comprises three key components: (1) Feedback signal identification, which detects and classifies user feedback, distinguishing between positive and negative signals to infer user preferences; (2) Preference data construction, which transforms these signals into structured preference datasets; and (3) Checklist-guided evaluation, which systematically assesses model responses using an instance-level checklist derived from extracted user preferences as a rubric. This ensures that model improvements are grounded in real user expectations rather than predefined heuristics.

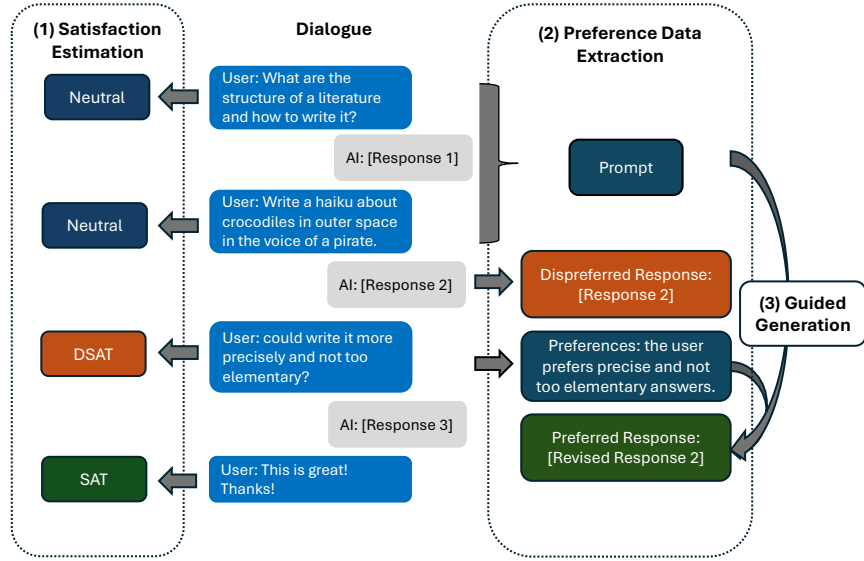


Figure 1: Overview of WILDFEEDBACK. (1) We begin by applying user satisfaction estimation to identify conversations and utterances that contain feedback signals. (2) We extract the entire conversation history leading up to a DSAT (dissatisfaction) signal as the prompt, and the response that triggers the DSAT as the dispreferred response. (3) Finally, we summarize the user’s preferences based on the identified feedback signals and guide the generation of the preferred response

To demonstrate the effectiveness of WILDFEEDBACK, we apply it to WildChat (Zhao et al., 2024), a dataset containing over 148,000 multi-turn conversations between users and ChatGPT (OpenAI et al., 2024) (see details of WildChat in Appendix E). This process results in a preference dataset of 20,281 samples¹, providing a rich resource for improving LLM alignment with real-world user preferences.

Through extensive experiments, we demonstrate that models fine-tuned on WILDFEEDBACK show significant improvements in aligning with user preferences, both in automated benchmarks and in our proposed checklist-guided evaluation framework. This work represents a step forward in creating more user-centric LLMs, with the potential to enhance user satisfaction across a wide range of applications.

The contributions of this paper are threefold:

- 1. Introduction of WILDFEEDBACK:** We present a novel framework that leverages in-situ user feedback to construct preference datasets that better reflect actual human values, addressing the scalability and subjectivity issues inherent in human-annotated datasets and the biases in synthetic data.

- 2. Robust Data Construction:** We adapt and ex-

pand on existing user satisfaction estimation techniques to identify feedback signals in natural conversations. This enables the creation of a nuanced preference dataset that includes both user preferences and corresponding responses, enhancing the effectiveness of fine-tuning LLMs to better align with user expectations.

- 3. Checklist-Guided Evaluation:** We propose a checklist-guided evaluation methodology that aligns the assessment of model performance with real user preferences, providing a more accurate benchmark for evaluating LLMs’ alignment with human values.

2 Related Work

Feedback Learning for LLMs. Incorporating human feedback has been shown to be an effective strategy to align LLMs with human preferences (Ouyang et al., 2022; Bai et al., 2022a; Dubey et al., 2024). However, relying human annotators to provide human feedback is inefficient and resource-intensive, which makes it hard to scale up. Additionally, human preferences are highly subjective. A small set of annotators may not represent broader preferences. Accordingly, some researchers aim to supervise AI models by model themselves (Bai et al., 2022b; Lee et al., 2023; Madaan et al., 2023; Burns et al., 2023; Li et al., 2023a). For instance,

¹The dataset will be available soon.

Bai et al. (2022b) introduced constitutional AI, in which they prompt LLMs to self-refine their own generations given a set of human-defined constitutions. However, relying on model’s own feedback can create a feedback loop where the model’s outputs increasingly reflect its own biases rather than diverse and authentic human perspectives. Recently, researchers have begun exploring the mining of user preferences from natural human-LLM interactions (Shi et al., 2022; Lin et al., 2024b; Don-Yehiya et al., 2024). These approaches capture real-time user feedback for more accurate preference alignment. Our work builds on this trend by leveraging in-situ user interactions to create preference datasets that better align with actual human values, addressing the limitations of both synthetic and human-annotated preference datasets.

Data for LLM Alignment. LLM alignment typically consists of two steps: instruction tuning and preference training. Instruction tuning, or supervised finetuning (SFT), aims to finetune models with a set of instruction-response pairs. Early works incorporated various NLP tasks for instruction tuning, demonstrating that LLMs could generalize well across different tasks (Wang et al., 2022; Chung et al., 2022; Ouyang et al., 2022). Subsequent research focused on constructing instruction data by directly distilling from capable LLMs (Wang et al., 2023; Xu et al., 2023). Researchers later recognized that preference training could further boost model performance across various tasks (Ouyang et al., 2022; Dubey et al., 2024). Preference training uses desired and undesired responses, either human-annotated (Bai et al., 2022a) or LLM-generated (Cui et al., 2024). Beyond general-purpose preference datasets, some datasets focus on specific tasks, such as summarization (Wu et al., 2021), model safety (Ji et al., 2023; Shi et al., 2024), and mathematics (Lightman et al., 2023). However, these approaches often rely on curated datasets that are either manually annotated by human experts or generated by models like GPT-4 (OpenAI et al., 2024). While these datasets provide a useful foundation, they may not fully capture the complexity and diversity of real-world user interactions. Our work addresses this gap by introducing a framework that leverages real-time feedback from actual users, allowing for more authentic and context-sensitive alignment of LLMs with true human preferences.

3 WILDFEEDBACK

Existing preference datasets often suffer from a mismatch between actual human preferences and those of the annotators (Chen et al., 2024; Poddar et al., 2024). Synthetic preference datasets, such as ULTRAFEEDBACK (Cui et al., 2024), rely solely on GPT-4 to generate rankings and determine which responses are preferred or dispreferred. However, this approach may not accurately capture real human values or nuanced preferences. Relying on synthetic data can create a feedback loop where the model’s outputs increasingly reflect its own biases rather than diverse and authentic human perspectives. On the other hand, preference datasets annotated by human annotators are difficult to scale due to time and budget constraints (Bai et al., 2022a; Ouyang et al., 2022; Dubey et al., 2024). Moreover, human annotators’ preferences can be highly subjective, often differing significantly from those of real users (Zhang et al., 2024; Fleisig et al., 2023).

To address these challenges, we introduce WILDFEEDBACK, a framework designed to align LLMs with in-situ user interactions and feedback. Unlike previous approaches that rely on synthetic responses, our framework directly learns preferences from real-world users, capturing both explicit and implicit feedback signals. The framework comprises three steps: (1) feedback signal identification, (2) preference data construction, and (3) checklist-guided evaluation. The pipeline is illustrated in Figure 1. We apply this framework to WildChat (Zhao et al., 2024), a corpus of real user-ChatGPT conversations, and obtained the WILDFEEDBACK dataset, a preference dataset of 20,281 samples.

3.1 Feedback Signals Identification

To construct preference data from natural human-LLM interactions, we first identify conversations that contain feedback signals. This can be achieved through user satisfaction estimation. In multi-turn conversational sessions, a user may explicitly express their satisfaction (e.g., “thank you”) or dissatisfaction (e.g., “revise it”) in their utterances. Lin et al. (2024b) proposed a framework named SPUR that can automatically learn and identify SAT (satisfaction) and DSAT (dissatisfaction) patterns. SPUR generalizes SAT/DSAT rubrics from conversations with annotated thumb feedback by recursively prompting GPT-4. These rubrics can then be used to score a user’s overall satisfaction or

Category	SAT	DSAT	Total
# Conversations	5,447	13,582	148,715
# Utterances	8,186	27,711	628,467

Table 1: Statistics of SAT/DSAT in conversations. A conversation is labeled as SAT/DSAT if it contains at least one SAT/DSAT utterance.

dissatisfaction, allowing us to identify utterances containing feedback signals.

WILD FEEDBACK adapts the SAT/DSAT rubrics from Lin et al. (2024b) with minor modifications. In total, we use 9 SAT and 9 DSAT rubrics. The SAT criteria include gratitude, learning, compliance, praise, personal details, humor, acknowledgment, positive closure, and getting there. The DSAT criteria consist of negative feedback, revision, factual error, unrealistic expectation, no engagement, ignored, lower quality, insufficient detail, and style. Detailed definitions of these rubrics can be found in Table 4 and Table 5. To streamline the process, we input these rubrics into GPT-4² and prompt it to perform the classification at the utterance level. The complete prompt is available in the Appendix A.1. In total, there are 148,715 multi-turn conversations in the WildChat dataset, with approximately 12.8% of the multi-turn conversations containing feedback signals. Detailed statistics are presented in Table 1.

To ensure the reliability of GPT-4’s classification of SAT/DSAT signals, we conducted a validation process using human expert annotators. Our findings indicate that GPT-4’s ability to identify SAT/DSAT signals shows relatively high agreement with human annotations, achieving a Cohen’s Kappa of $\kappa = 0.69$ for SAT and $\kappa = 0.50$ for DSAT, similar to the human performance. A detailed breakdown of GPT-4’s performance and the human annotation process are provided in Appendix B.2.

3.2 Preference Pair Generation

After identifying conversations that contain feedback signals using the SAT/DSAT rubrics, we can construct semi-synthetic preference pairs. Each preference pair sample consists of four components: the prompt, user preferences, the preferred response, and the dispreferred response. For con-

versations with SAT/DSAT signals, we first analyze user responses marked by these signals and ask GPT-4 to summarize user preferences based on these feedback signals (e.g., the user prefers concise and direct answers). We then extract the conversation up to the model response that triggers the SAT/DSAT signals and use this as the prompt for our preference data.

For preferred and dispreferred response generation, we explore two different approaches: expert responses and on-policy responses. Specifically, we use GPT-4 for expert response generation, while Phi 3 (Abdin et al., 2024), Qwen 2 (Yang et al., 2024), and LLaMA 3 (Dubey et al., 2024) are employed for on-policy response generation. For expert responses, those that trigger DSAT signals in the original conversations are directly used as dispreferred responses (e.g., response 2 in Fig. 1). We then prompt GPT-4 to generate the preferred responses by using summarized user preferences as the system prompt. For on-policy responses, both preferred and dispreferred responses are generated by the policy model. The dispreferred responses are generated directly, whereas the preferred responses are produced using the summarized user preferences as the system prompt. Furthermore, recognizing that some user preferences may be harmful (e.g., preferences for explicit content), we take extra safety precautions. When prompting either the on-policy models or GPT-4 to generate preferred responses, we include an additional system instruction: “The response should be safe.” Some conversations are also automatically filtered by the OpenAI moderation API. The prompt used for preference pair construction is provided in Appendix A.2.

3.3 Checklist-guided Evaluation

Existing automated benchmarks, such as AlpacaEval (Dubois et al., 2024) and MT-Bench (Zheng et al., 2023b), heavily rely on using LLMs as judges. These benchmarks typically prompt models with a set of queries and then ask LLMs like GPT-4 or Claude (Anthropic, 2023) to provide a score or rank the responses of different models. This approach is problematic because it relies heavily on the internal knowledge of LLMs, which are known to be biased towards longer responses or responses generated by themselves (Liu et al., 2024b; Thakur et al., 2024). Additionally, there is a mismatch between the preferences of LLMs as judges and those of humans, leading to evaluations that

²Unless otherwise specified, in all of our experiments, we use GPT-4o with the gpt-4o-0513 engine. For open-weight models, we use Phi-3-mini-4k-instruct, Qwen2-7B-Instruct, Meta-Llama-3-8B-Instruct.

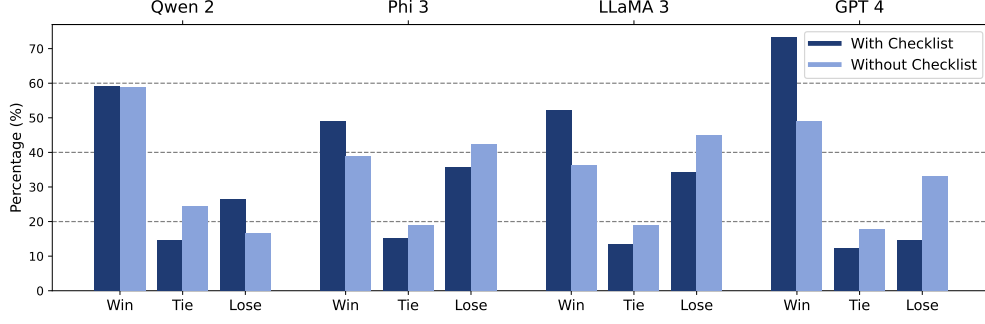


Figure 2: Comparison of in-situ user alignment across datasets generated by different models. “Win/Tie/Lose” represents the percentage of instances where the preferred responses win/tie/lose compared to the dispreferred responses in the WILDfeedback dataset, prior to filtering. The comparison is made both with and without providing GPT-4 with summarized user preferences as checklists to guide its evaluation. With checklists, the preferred responses can be better distinguished.

do not accurately reflect user preferences. Furthermore, using human annotators to rank model responses based on their subjective experiences is also not ideal, as there can be a mismatch between annotators’ preferences and actual user preferences.

In response, we propose checklist-guided evaluation, a general evaluation framework that more accurately reflects real user preferences. In our preference data construction module, we not only construct preference data from user-LLM interactions but also summarize user preferences expressed in natural language. These preferences, based on real users’ textual feedback, can be used to align LLMs’s evaluation more closely with real users’ preferences. Instead of asking human annotators to directly rank model responses, we should ask them to rank those responses based on real users’ preferences. When using LLMs as evaluators, we can provide an instance-level checklist to guide their assessments. Our evaluation framework is adapted from WILDBENCH (Lin et al., 2024a), which has been shown to correlate well with human judgement in ranking model performance as an automatic metric. We employ a pairwise evaluation strategy, where GPT-4 compares two different responses to determine which performs better on a given task, using an instance-level, preference-guided checklist to inform the comparison. This metric allows for straightforward comparisons among models, with easily interpretable win/lose rates as intermediate outcomes. The full prompt can be found in Appendix A.3.

Similar to feedback signal identification (§3.1), to ensure the reliability of GPT-4 on checklist-guided evaluation, we conducted a validation pro-

cess using human expert annotators. We found GPT-4 achieves an human agreement of 57.14%, similar to the human-human agreement of 63.27%. A detailed breakdown of GPT-4’s performance and the human annotation process are provided in Appendix C.

3.4 WILDfeedback Data Construction

The preference pair construction approach described in Section 3.2 allows us to build a robust dataset for training models to better align responses with user preferences.

To evaluate whether our generated preferred responses align with actual user preferences, we randomly selected 500 samples from the WILDfeedback datasets and performed checklist-guided evaluation (§3.3), comparing the preferred and dispreferred responses. As explained in Section 3.2, there are two versions of WILDfeedback preference pairs: the GPT-4 version and the on-policy version, which differ in whether the responses are generated by GPT-4 or the policy model. As shown in Figure 2, we found that without checklist-guided evaluation, GPT-4 does not necessarily favor responses aligned with summarized user preferences, often defaulting to models’ zero-shot generations instead. However, after providing the preferences as checklists to guide the evaluation, GPT-4’s selections more closely align with real users’ preferences. Additionally, we observed that GPT-4 is significantly more steerable than smaller models: over 70% of its preferred responses align with in-situ user preferences, compared to only about 50% for smaller models.

Since policy models are less steerable than GPT-4 and may not always align with provided user

	# Conv.	Prompt Length	Response Length	Multi-Turn?	Feedback Type
WebGPT (Nakano et al., 2022)	38,925	51	188	✗	Human Annotators
Anthropic HH (Bai et al., 2022a)	118,263	186	95	✗	Human Annotators
OASST1 (Köpf et al., 2023)	35,905	168	221	✓	Human Annotators
ULTRAFEEDBACK (Cui et al., 2024)	61,135	159	256	✗	GPT-4
WILDFEEDBACK (ours)					
↪ GPT-4	20,281	929	440		
↪ Qwen 2	11,509	1,057	541	✓	In-situ Users
↪ Phi 3	9,194	931	344		
↪ LLaMA 3	10,659	982	376		

Table 2: Statistics of existing preference datasets. The average length refers to the number of tokens. The responses of WILDFEEDBACK are either extracted from the original conversations or generated by GPT-4, Qwen 2, Phi 3, or LLaMA 3.

preferences, we apply an additional filtering process, discarding any on-policy pairs that do not align with user preferences based on checklist-guided evaluation. In contrast, we retain all GPT-4-generated preference pairs, as they consistently demonstrate higher alignment.

Table 2 reports statistics on WILDFEEDBACK constructed datasets compared with open-source datasets³. To the best of our knowledge, WILDFEEDBACK is the first multi-turn pairwise preference dataset derived from real human-LLM interactions. Unlike datasets annotated by human labelers or LLMs, which often fail to fully capture real user preferences, WILDFEEDBACK is built from in-situ user feedback. Although OpenAssistant Conversations (OASST1) (Köpf et al., 2023) also includes multi-turn conversations, its prompts and responses are fully composed by human annotators, making it less reflective of genuine human-LLM interactions. In the next section, we demonstrate that WILDFEEDBACK more accurately represents authentic human-LLM interactions, making it a more reliable resource for developing and evaluating preference-based models.

4 Experiment

To validate the effectiveness of WILDFEEDBACK, we finetune models from different families on it and compare their performances with the vanilla models and the models finetuned on ULTRAFEEDBACK data. We evaluate models’ performance on general benchmarks and a held-out test set of WILDFEEDBACK using checklist-guided evaluation.

Models and training settings. We use off-the-shelf instruction-tuned Qwen 2, Phi 3, and LLaMA

3 models. As described in Section 3.2, each model is fine-tuned on two versions of both WILDFEEDBACK (WF) and ULTRAFEEDBACK (UF): a GPT-4 version and an on-policy version.

For WILDFEEDBACK, the WF GPT-4 setup utilizes GPT-4 to generate preferred responses based on summarized user preferences. Dispreferred responses are extracted from conversations that contain DSAT signals. In the WF On-policy setup, each policy model (Qwen 2, Phi 3, or LLaMA 3) generates both preferred and dispreferred responses, again making use of summarized user preferences to produce the preferred ones. We train each model for one epoch of supervised fine-tuning (SFT) on the preferred responses, followed by one epoch of direct preference optimization (DPO) (Rafailov et al., 2023) on the entire dataset. We find that hyperparameter tuning is essential for optimal results (see Appendix D).

We also fine-tune models using ULTRAFEEDBACK, one of the most widely used preference datasets due to its superior performance compared to others. Models such as the Tulu 3 series (Lambert et al., 2025) and Zephyr (Tunstall et al., 2023) have been fine-tuned on this dataset. The prompts in ULTRAFEEDBACK are sourced from various instruction datasets. Each prompt has four responses from different LLMs, numerically rated by GPT-4. However, due to the off-policy nature of ULTRAFEEDBACK and the outdated models used to generate its responses, it has become common practice to regenerate responses using only the original prompts when training new models on this dataset (Meng et al., 2024; Dong et al., 2024; Xiong et al., 2024). Following this approach, we create two versions of the dataset: UF GPT-4 and UF On-policy. In UF GPT-4, we randomly select 20,000 prompts from ULTRAFEEDBACK, and GPT-4 generates two

³For ULTRAFEEDBACK, we refer to the pre-processed, binarized version used to train Zephyr (Tunstall et al., 2023).

responses for each prompt. GPT-4 then acts as a judge, selecting the better response as the preferred one while marking the other as dispreferred. In UF On-policy, each policy model generates five responses per prompt, after which a GPT-4 judge selects the best response as preferred, while one of the remaining four is randomly designated as dispreferred. The specific prompt used to guide GPT-4 in selecting the preferred response is provided in Appendix A.4. By regenerating the responses for ULTRAFEEDBACK, we also ensure a fair comparison to our WILDFEEDBACK setup.

In summary, for all three policy models, we compare five configurations: (1) the off-the-shelf instruction-tuned model, (2) WF GPT-4, (3) WF On-policy, (4) UF GPT-4, and (5) UF On-policy.

Benchmarks Evaluation. We evaluate our models using three of the most popular open-ended instruction-following benchmarks: AlpacaEval 2 (Li et al., 2023b), MT-Bench (Zheng et al., 2023a), and Arena-Hard (Li et al., 2024). AlpacaEval 2 consists of 805 questions from 5 datasets, and MT-Bench covers 8 categories with 80 questions. Arena-Hard is an enhanced version of MT-Bench, incorporating 500 well-defined technical problem-solving queries. We report scores following each benchmark’s evaluation protocol: For AlpacaEval 2, we report both the raw win rate (WR) and the length-controlled win rate (LC) (Dubois et al., 2024). The LC metric is specifically designed to be robust against model verbosity. For MT-Bench, we report the average MT-Bench score with GPT-4o (gpt-4o-0513) as the judge. For Arena-Hard, we report the win rate (WR) against the baseline model. As specified by the benchmarks, we use GPT-4-Turbo (gpt-4-0125) as the judge for both AlpacaEval 2 and Arena-Hard. We use the same, default decoding strategy specified by each evaluation benchmark respectively.

WILDFEEDBACK Evaluation. In addition to publicly available benchmarks, we constructed our own evaluation benchmark from the held-out test set in WILDFEEDBACK and evaluated models using checklist-guided evaluation (§3.3). We ensured that all test samples came from conversations and users that were never included in the training set. Constructing an evaluation dataset for checklist-guided evaluation is non-trivial, as we can no longer randomly or stratifiedly select test samples from different domains. In checklist-guided evaluation, we always provide a user-inspired checklist

for GPT-4 to guide its evaluation, making it more aligned with real users’ preferences. However, individual user preferences can be highly subjective and specific. The goal of WILDFEEDBACK is not to align language models with the preferences of a specific individual but to learn the broader mode of all individuals’ preferences. Therefore, we must ensure that the preferences reflected in the test samples represent the majority view. Additionally, since the user preferences we extracted are often particular to specific tasks, we also need to ensure that the tasks in the test set are at least somewhat similar to those in the training set.

To achieve this, we utilized FAISS (Douze et al., 2024) to cluster user prompts and their summarized preferences. We grouped all user prompts into 70 clusters. Within each cluster, we selected 10 samples where the preferences were most similar to the other preferences in the same group. We then applied similar data curation techniques as described in WILDBENCH (Lin et al., 2024a) to perform deduplication and remove nonsensical tasks, resulting in a final test set of 540 samples. By doing so, we aim to provide a more reliable and comprehensive evaluation that reflects the majority’s preferences without overfitting to specific, idiosyncratic cases.

For WILDFEEDBACK evaluation, we report the win, tie, lose percentage against the instruct models and the models trained on ULTRAFEEDBACK with GPT-4 as the judge. We employ the WILDBENCH prompt (Lin et al., 2024a) to perform the evaluation, which has been shown to correlate well with human judgement in ranking model performance. We report the results evaluated with or without the user preferences provided as a checklist.

5 Results and Analysis

In this section, we present the main results of our experiments, highlighting the effectiveness of WILDFEEDBACK on various benchmarks and ablation studies.

Training models on WILDFEEDBACK significantly and consistently enhances performance across all benchmarks. As shown in Table 3, models trained on either version of WILDFEEDBACK achieve higher performance across AlpacaEval 2, Arena-Hard, and MT-Bench. For example, after training on the GPT-4 version of WILDFEEDBACK (WF GPT-4), Phi 3’s length-controlled win rate on AlpacaEval 2 increases from 24.3% to

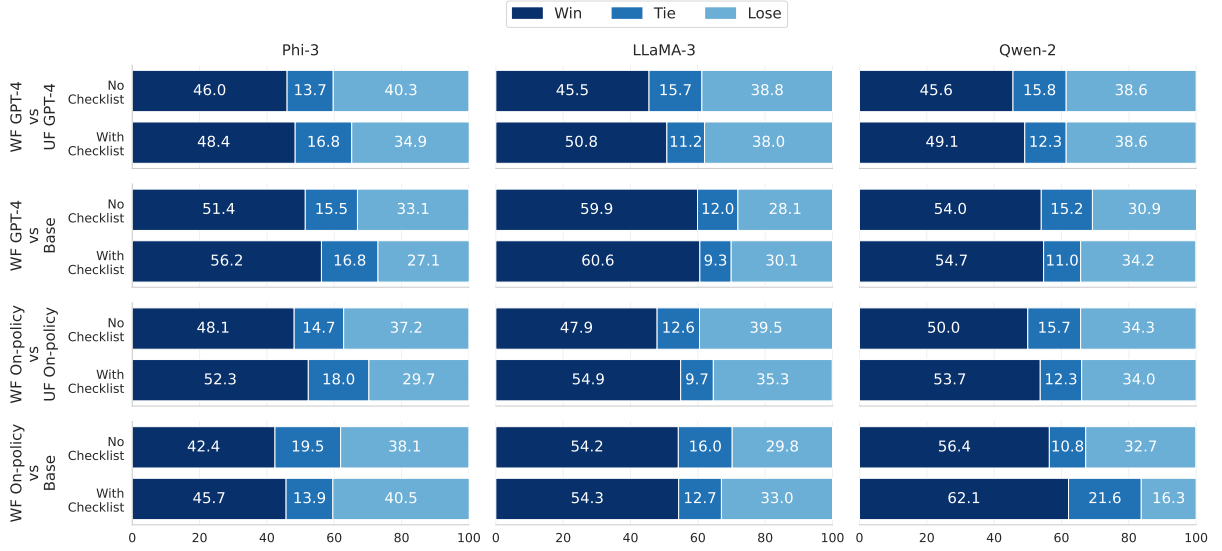


Figure 3: Preference evaluation on the WILD FEEDBACK test set, with or without the checklist. All numbers are the percentages of win/tie/lose. WF/UF On-policy/GPT-4 refers to the model trained on the on-policy/GPT-4 version of WILD FEEDBACK/ULTRA FEEDBACK. Base models here refers to the off-the-shelf instruct models. Models trained on WILD FEEDBACK consistently outperformed all the baselines.

34.9%, while its win rate on Arena-Hard improves from 15.4% to 32.4%. Similarly, its performance on MT-Bench rises from a score of 7.32 to 7.75. Models trained on WILD FEEDBACK also consistently outperform those on ULTRA FEEDBACK.

Models	AlpacaEval 2		Arena-Hard	MT-Bench
	LC (%)	WR (%)	WR (%)	Score
Phi 3	24.3	17.4	15.4	7.32
↪ WF On-Policy	29.0	27.1	30.1	7.42
↪ UF On-Policy	27.2	25.9	28.7	7.40
↪ WF GPT-4	34.9	36.6	32.4	7.75
↪ UF GPT-4	32.5	38.4	30.5	7.68
LLaMA 3	22.9	22.6	20.6	7.10
↪ WF On-Policy	30.1	29.6	22.1	7.15
↪ UF On-Policy	28.8	34.1	20.2	7.04
↪ WF GPT-4	34.2	42.8	32.9	7.57
↪ UF GPT-4	32.2	43.2	32.6	7.49
Qwen 2	28.7	26.0	24.9	7.55
↪ WF On-Policy	42.6	34.4	36.1	8.02
↪ UF On-Policy	38.3	34.2	29.2	7.72
↪ WF GPT-4	39.4	33.5	27.9	7.60
↪ UF GPT-4	40.6	32.5	27.6	7.66

Table 3: AlpacaEval 2, Arena-Hard, and MT-Bench results under the four settings. LC and WR denote length-controlled and raw win rate. WF/UF On-policy/GPT-4 refers to the model trained on the on-policy/GPT-4 version of WILD FEEDBACK/ULTRA FEEDBACK.

WILD FEEDBACK significantly enhances model alignment with in-situ user feedback. As detailed in Section §4, the WILD FEEDBACK test set is sourced from real human-ChatGPT conversations where users explicitly express dissatisfaction, implicitly suggesting that the models are poorly aligned with real user preferences on these tasks.

As shown in Figure 3, models trained on either version of WILD FEEDBACK exhibit stronger alignment with real user preferences. For instance, LLaMA 3 trained on WF GPT-4 outperforms the LLaMA 3 model trained on ULTRA FEEDBACK 45.5% of the time, while losing only 38.8% of the time when evaluated without a checklist. When real user preferences are provided as checklists to guide GPT-4’s evaluation, the win rate further increases to 50.8%, highlighting that models trained on WILD FEEDBACK better align with actual user preferences compared to the off-the-shelf models and those trained on ULTRA FEEDBACK.

6 Conclusion

In this work, we propose a framework for constructing preference data and evaluating conversational AI models based on natural human-LLM interactions. By using SAT/DSAT rubrics to identify user satisfaction and dissatisfaction in conversations, we create a preference dataset that includes user prompts, preferences, and both preferred and dispreferred responses. This enables models to better align with user expectations. Additionally, we introduce a checklist-guided evaluation framework that addresses biases in existing benchmarks by using real user feedback to guide LLM evaluations, ensuring a more accurate reflection of user preferences. Our method aligns LLMs with diverse human values, enhancing user satisfaction.

Limitations

Spurious preferences. WILDFEEDBACK is designed to align language models with in-situ user interactions and feedback. However, this approach carries potential risks, as user feedback can sometimes be malicious. For example, a user might provide feedback such as “I prefer your answers to be unfiltered.” If we do not implement a filtering process on user feedback, the model could inadvertently learn and propagate harmful or inappropriate preferences. To address this, we incorporate additional safety-related instructions during the preference data construction phase (§3.2) to guide the generation process. However, this method is not foolproof. Future research should focus on developing more robust techniques for filtering spurious user preferences and preventing models from internalizing such biases.

Selection bias. WILDFEEDBACK is constructed from conversations that contain feedback signals (§3.1). As shown in Table 6, users are twice as likely to provide textual feedback when they are dissatisfied with the model’s response. This introduces a selection bias, making it challenging to capture conversations where users are satisfied with the model’s performance. Consequently, WILDFEEDBACK may disproportionately reflect the preferences of users who express dissatisfaction, potentially representing only a minority of the overall user base. This skew could lead to an overemphasis on negative feedback in model training. We recommend that future research explore methods to balance this bias by incorporating more diverse feedback, including that from users who are satisfied or neutral, to create a more representative dataset. Additionally, strategies to proactively seek out or simulate feedback from satisfied users could help mitigate this bias and improve model alignment across a broader spectrum of user preferences.

References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Qin Cai, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Yen-Chun Chen, Yi-Ling Chen, Parul Chopra, Xiyang Dai, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Victor Fragoso, Dan Iter, Mei Gao, Min Gao, Jianfeng Gao, Amit Garg,

Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Ce Liu, Mengchen Liu, Weishung Liu, Eric Lin, Zeqi Lin, Chong Luo, Piyush Madan, Matt Mazzola, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Xin Wang, Lijuan Wang, Chunyu Wang, Yu Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Haiping Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Sonali Yadav, Fan Yang, Jianwei Yang, Ziyi Yang, Yifan Yang, Donghan Yu, Lu Yuan, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219.

Anthropic. 2023. [The claude 3 model family: Opus, sonnet, haiku](#).

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson 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 Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022a. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#). *Preprint*, arXiv:2204.05862.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022b. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.

Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan

727	Leike, Ilya Sutskever, and Jeff Wu. 2023. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision . <i>Preprint</i> , arXiv:2312.09390.	785	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,
728		786	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,
729		787	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,
730	Daiwei Chen, Yi Chen, Aniket Rege, and Ramya Korlakai Vinayak. 2024. Pal: Pluralistic alignment framework for learning from heterogeneous preferences . <i>Preprint</i> , arXiv:2406.08469.	788	Emily Dinan, Eric Michael Smith, Filip Radenovic,
731		789	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, 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, 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 Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier 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, Roman Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Conguet, 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 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
732		790	
733		791	
734	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models . <i>Preprint</i> , arXiv:2210.11416.	792	
735		793	
736		794	
737		795	
738		796	
739		797	
740		798	
741		799	
742		800	
743		801	
744		802	
745		803	
746		804	
747	Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. ULTRA FEEDBACK: Boosting language models with scaled AI feedback . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 9722–9744. PMLR.	805	
748		806	
749		807	
750		808	
751		809	
752		810	
753		811	
754	Sarkar Snigdha Sarathi Das, Chirag Shah, Mengting Wan, Jennifer Neville, Longqi Yang, Reid Andersen, Georg Buscher, and Tara Safavi. 2023. S3dst: Structured open-domain dialogue segmentation and state tracking in the era of llms . <i>Preprint</i> , arXiv:2309.08827.	812	
755		813	
756		814	
757		815	
758		816	
759		817	
760	Shachar Don-Yehiya, Leshem Choshen, and Omri Abend. 2024. Learning from naturally occurring feedback . <i>Preprint</i> , arXiv:2407.10944.	818	
761		819	
762		820	
763	Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. 2024. RLhf workflow: From reward modeling to online rlhf . <i>Preprint</i> , arXiv:2405.07863.	821	
764		822	
765		823	
766		824	
767		825	
768	Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library . <i>Preprint</i> , arXiv:2401.08281.	826	
769		827	
770		828	
771		829	
772	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 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 Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	830	
773		831	
774		832	
775		833	
776		834	
777		835	
778		836	
779		837	
780		838	
781		839	
782		840	
783		841	
784		842	
		843	
		844	
		845	
		846	
		847	
		848	

849	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, Changan 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, Dingkan 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 Flores, Gabriella Schwarz, Gada Badeer, Georgia Sweet, 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, 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 Michelen, 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 Khabza, Manav Avalani, Manish Bhatt, Maria Tsim-poukelli, 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 Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratan-chandani, 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	
850		
851		
852		
853		
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		
864		
865		
866		
867		
868		
869		
870		
871		
872		
873		
874		
875		
876		
877		
878		
879		
880		
881		
882		
883		
884		
885		
886		
887		
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
	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, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan 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 . <i>Preprint</i> , arXiv:2407.21783.	913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937
	Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators . <i>Preprint</i> , arXiv:2404.04475.	938 939 940 941
	Eve Fleisig, Rediet Abebe, and Dan Klein. 2023. When the majority is wrong: Modeling annotator disagreement for subjective tasks . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6715–6726, Singapore. Association for Computational Linguistics.	942 943 944 945 946 947
	Sanjana Gautam and Mukund Srinath. 2024. Blind spots and biases: Exploring the role of annotator cognitive biases in nlp . <i>Preprint</i> , arXiv:2404.19071.	948 949 950
	Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Chi Zhang, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. Beavertails: Towards improved safety alignment of llm via a human-preference dataset . <i>Preprint</i> , arXiv:2307.04657.	951 952 953 954 955
	Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Minh Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Alexandrovich Glushkov, Arnav Varma Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Julian Mattick. 2023. Openassistant conversations - democratizing large language model alignment . In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	956 957 958 959 960 961 962 963 964 965 966
	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le	967 968 969 970 971

972	Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini,	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	1028
973	Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	1029
974	Hannaneh Hajishirzi. 2025. Tulu 3: Pushing fron-	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	1030
975	tiers in open language model post-training . <i>Preprint</i> ,	Shashank Gupta, Bodhisattwa Prasad Majumder,	1031
976	arXiv:2411.15124.	Katherine Hermann, Sean Welleck, Amir Yazdan-	1032
977	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas	bakhsh, and Peter Clark. 2023. Self-refine: It-	1033
978	Mesnard, Johan Ferret, Kellie Lu, Colton Bishop,	erative refinement with self-feedback . <i>Preprint</i> ,	1034
979	Ethan Hall, Victor Carbune, Abhinav Rastogi, and	arXiv:2303.17651.	1035
980	Sushant Prakash. 2023. Rlaif: Scaling reinforce-	Yu Meng, Mengzhou Xia, and Danqi Chen. 2024.	1036
981	ment learning from human feedback with ai feedback .	Simpo: Simple preference optimization with a	1037
982	<i>Preprint</i> , arXiv:2309.00267.	reference-free reward . <i>Preprint</i> , arXiv:2405.14734.	1038
983	Minzhi Li, Taiwei Shi, Caleb Ziems, Min-Yen Kan,	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu,	1039
984	Nancy Chen, Zhengyuan Liu, and Diyi Yang. 2023a.	Long Ouyang, Christina Kim, Christopher Hesse,	1040
985	CoAnnotating: Uncertainty-guided work allocation	Shantanu Jain, Vineet Kosaraju, William Saunders,	1041
986	between human and large language models for data	Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen	1042
987	annotation . In <i>Proceedings of the 2023 Conference</i>	Krueger, Kevin Button, Matthew Knight, Benjamin	1043
988	<i>on Empirical Methods in Natural Language Process-</i>	Chess, and John Schulman. 2022. Webgpt: Browser-	1044
989	<i>ing</i> , pages 1487–1505, Singapore. Association for	assisted question-answering with human feedback .	1045
990	Computational Linguistics.	<i>Preprint</i> , arXiv:2112.09332.	1046
991	Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap,	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	1047
992	Tianhao Wu, Banghua Zhu, Joseph E. Gonzalez, and	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	1048
993	Ion Stoica. 2024. From crowdsourced data to high-	man, Diogo Almeida, Janko Altmenschmidt, Sam Alt-	1049
994	quality benchmarks: Arena-hard and benchbuilder	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	1050
995	pipeline . <i>Preprint</i> , arXiv:2406.11939.	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	1051
996	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori,	ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-	1052
997	Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	1053
998	Tatsunori B. Hashimoto. 2023b. AlpacaEval: An	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	1054
999	automatic evaluator of instruction-following models .	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	1055
1000	https://github.com/tatsu-lab/alpaca_eval .	man, Tim Brooks, Miles Brundage, Kevin Button,	1056
1001	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	1057
1002	Edwards, Bowen Baker, Teddy Lee, Jan Leike,	Carey, Chelsea Carlson, Rory Carmichael, Brooke	1058
1003	John Schulman, Ilya Sutskever, and Karl Cobbe.	Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	1059
1004	2023. Let’s verify step by step . <i>Preprint</i> ,	Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	1060
1005	arXiv:2305.20050.	Chess, Chester Cho, Casey Chu, Hyung Won Chung,	1061
1006	Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze	Dave Cummings, Jeremiah Currier, Yunxing Dai,	1062
1007	Brahman, Abhilasha Ravichander, Valentina Py-	Cory Decareaux, Thomas Degry, Noah Deutsch,	1063
1008	atkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi.	Damien Deville, Arka Dhar, David Dohan, Steve	1064
1009	2024a. Wildbench: Benchmarking llms with chal-	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	1065
1010	lenging tasks from real users in the wild . <i>Preprint</i> ,	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	1066
1011	arXiv:2406.04770.	Simón Posada Fishman, Juston Forte, Isabella Ful-	1067
1012	Ying-Chun Lin, Jennifer Neville, Jack W. Stokes,	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	1068
1013	Longqi Yang, Tara Safavi, Mengting Wan, Scott	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	1069
1014	Counts, Siddharth Suri, Reid Andersen, Xiaofeng	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	1070
1015	Xu, Deepak Gupta, Sujay Kumar Jauhar, Xia Song,	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	1071
1016	Georg Buscher, Saurabh Tiwary, Brent Hecht, and	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	1072
1017	Jaime Teevan. 2024b. Interpretable user satisfaction	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	1073
1018	estimation for conversational systems with large lan-	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	1074
1019	guage models . <i>Preprint</i> , arXiv:2403.12388.	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	1075
1020	Jie Liu, Zhanhui Zhou, Jiaheng Liu, Xingyuan Bu, Chao	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	1076
1021	Yang, Han-Sen Zhong, and Wanli Ouyang. 2024a. It-	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	1077
1022	erative length-regularized direct preference optimiza-	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	1078
1023	tion: A case study on improving 7b language models	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	1079
1024	to gpt-4 level . <i>Preprint</i> , arXiv:2406.11817.	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	1080
1025	Yiqi Liu, Nafise Sadat Moosavi, and Chenghua Lin.	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	1081
1026	2024b. Llms as narcissistic evaluators: When ego in-	Christina Kim, Yongjik Kim, Jan Hendrik Kirch-	1082
1027	flates evaluation scores . <i>Preprint</i> , arXiv:2311.09766.	ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo,	1083
		Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	1084
		stantinidis, Kyle Kosic, Gretchen Krueger, Vishal	1085
		Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	1086
		Leike, Jade Leung, Daniel Levy, Chak Ming Li,	1087
		Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	1088

1089	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	portable text annotation tool. In <i>Proceedings of the</i>	1150
1090	Anna Makanju, Kim Malfacini, Sam Manning, Todor	<i>2022 Conference on Empirical Methods in Natural</i>	1151
1091	Markov, Yaniv Markovski, Bianca Martin, Katie	<i>Language Processing: System Demonstrations.</i>	1152
1092	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer		
1093	McKinney, Christine McLeavey, Paul McMillan,	Sriyash Poddar, Yanming Wan, Hamish Ivison, Ab-	1153
1094	Jake McNeil, David Medina, Aalok Mehta, Jacob	hishek Gupta, and Natasha Jaques. 2024. Person-	1154
1095	Menick, Luke Metz, Andrey Mishchenko, Pamela	alizing reinforcement learning from human feedback	1155
1096	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	with variational preference learning. In <i>Pluralistic</i>	1156
1097	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	<i>Alignment Workshop at NeurIPS 2024.</i>	1157
1098	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,		
1099	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	1158
1100	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	pher D Manning, Stefano Ermon, and Chelsea Finn.	1159
1101	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	2023. Direct preference optimization: Your language	1160
1102	tista Parascandolo, Joel Parish, Emy Parparita, Alex	model is secretly a reward model. In <i>Thirty-seventh</i>	1161
1103	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	<i>Conference on Neural Information Processing Sys-</i>	1162
1104	man, Filipe de Avila Belbute Peres, Michael Petrov,	<i>tems.</i>	1163
1105	Henrique Ponde de Oliveira Pinto, Michael, Poko-		
1106	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	Taiwei Shi, Kai Chen, and Jieyu Zhao. 2024. Safer-	1164
1107	ell, Alethea Power, Boris Power, Elizabeth Proehl,	instruct: Aligning language models with automated	1165
1108	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	preference data. In <i>Proceedings of the 2024 Confer-</i>	1166
1109	Cameron Raymond, Francis Real, Kendra Rimbach,	<i>ence of the North American Chapter of the Associ-</i>	1167
1110	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	<i>ation for Computational Linguistics: Human Lan-</i>	1168
1111	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	<i>guage Technologies (Volume 1: Long Papers),</i> pages	1169
1112	Girish Sastry, Heather Schmidt, David Schnurr, John	7636–7651, Mexico City, Mexico. Association for	1170
1113	Schulman, Daniel Selsam, Kyla Sheppard, Toki	Computational Linguistics.	1171
1114	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav		
1115	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	Weiyan Shi, Emily Dinan, Kurt Shuster, Jason We-	1172
1116	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	ston, and Jing Xu. 2022. When life gives you	1173
1117	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	lemons, make cherryade: Converting feedback	1174
1118	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	from bad responses into good labels. <i>Preprint,</i>	1175
1119	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	arXiv:2210.15893.	1176
1120	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,		
1121	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	Gemma Team, Morgane Riviere, Shreya Pathak,	1177
1122	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupati-	1178
1123	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	raju, Léonard Hussenot, Thomas Mesnard, Bobak	1179
1124	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	Shahriari, Alexandre Ramé, Johan Ferret, Peter	1180
1125	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	Liu, Pouya Tafti, Abe Friesen, Michelle Casbon,	1181
1126	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	Sabela Ramos, Ravin Kumar, Charline Le Lan,	1182
1127	Clemens Winter, Samuel Wolrich, Hannah Wong,	Sammy Jerome, Anton Tsitsulin, Nino Vieillard,	1183
1128	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	Piotr Stanczyk, Sertan Girgin, Nikola Momchev,	1184
1129	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill,	1185
1130	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	Behnam Neyshabur, Olivier Bachem, Alanna Wal-	1186
1131	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	ton, Aliaksei Severyn, Alicia Parrish, Aliya Ah-	1187
1132	Zheng, Juntang Zhuang, William Zhuk, and Bar-	mad, Allen Hutchison, Alvin Abdagic, Amanda	1188
1133	ret Zoph. 2024. Gpt-4 technical report. <i>Preprint,</i>	Carl, Amy Shen, Andy Brock, Andy Coenen, An-	1189
1134	arXiv:2303.08774.	thony Laforge, Antonia Paterson, Ben Bastian, Bilal	1190
		Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu	1191
1135	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	Kumar, Chris Perry, Chris Welty, Christopher A.	1192
1136	roll L. Wainwright, Pamela Mishkin, Chong Zhang,	Choquette-Choo, Danila Sinopalnikov, David Wein-	1193
1137	Sandhini Agarwal, Katarina Slama, Alex Ray, John	berger, Dimple Vijaykumar, Dominika Rogozińska,	1194
1138	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	Dustin Herbison, Elisa Bandy, Emma Wang, Eric	1195
1139	Maddie Simens, Amanda Askell, Peter Welinder,	Noland, Erica Moreira, Evan Senter, Evgenii Elty-	1196
1140	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	shev, Francesco Visin, Gabriel Rasskin, Gary Wei,	1197
1141	Training language models to follow instructions with	Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna	1198
1142	human feedback. <i>Preprint,</i> arXiv:2203.02155.	Klimczak-Plucińska, Harleen Batra, Harsh Dhand,	1199
		Ivan Nardini, Jacinda Mein, Jack Zhou, James Svens-	1200
1143	Jupinder Parmar, Sanjev Satheesh, Mostofa Patwary,	son, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana	1201
1144	Mohammad Shoeybi, and Bryan Catanzaro. 2024.	Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fer-	1202
1145	Reuse, don’t retrain: A recipe for continued pretrain-	nandez, Joost van Amersfoort, Josh Gordon, Josh	1203
1146	ing of language models. <i>Preprint,</i> arXiv:2407.07263.	Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mo-	1204
		hamed, Kartikeya Badola, Kat Black, Katie Mil-	1205
1147	Jiaxin Pei, Aparna Ananthasubramaniam, Xingyao	lican, Keelin McDonnell, Kelvin Nguyen, Kiranbir	1206
1148	Wang, Naitian Zhou, Apostolos Dedeloudis, Jack-	Sodhia, Kish Greene, Lars Lowe Sjoesund, Lau-	1207
1149	son Sargent, and David Jurgens. 2022. Potato: The	ren Usui, Laurent Sifre, Lena Heuermann, Leti-	1208
		cia Lago, Lilly McNealus, Livio Baldini Soares,	1209

1210	Logan Kilpatrick, Lucas Dixon, Luciano Martins,	Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma,	1271
1211	Machel Reid, Manvinder Singh, Mark Iverson, Mar-	Ravsehaj Singh Puri, Rushang Karia, Savan Doshi,	1272
1212	tin Görner, Mat Velloso, Mateo Wirth, Matt Davi-	Shailaja Keyur Sampat, Siddhartha Mishra, Sujan	1273
1213	dow, Matt Miller, Matthew Rahtz, Matthew Wat-	Reddy A, Sumanta Patro, Tanay Dixit, and Xudong	1274
1214	son, Meg Risdal, Mehran Kazemi, Michael Moyni-	Shen. 2022. Super-NaturalInstructions: Generaliza-	1275
1215	han, Ming Zhang, Minsuk Kahng, Minwoo Park,	tion via declarative instructions on 1600+ NLP tasks.	1276
1216	Mofi Rahman, Mohit Khatwani, Natalie Dao, Nen-	In <i>Proceedings of the 2022 Conference on Empiri-</i>	1277
1217	shad Bardoliwalla, Nesh Devanathan, Neta Dumai,	<i>cal Methods in Natural Language Processing</i> , pages	1278
1218	Nilay Chauhan, Oscar Wahltinez, Pankil Botarda,	5085–5109, Abu Dhabi, United Arab Emirates. As-	1279
1219	Parker Barnes, Paul Barham, Paul Michel, Peng-	sociation for Computational Linguistics.	1280
1220	chong Jin, Petko Georgiev, Phil Culliton, Pradeep		
1221	Kuppala, Ramona Comanescu, Ramona Merhej,	Jeff Wu, Long Ouyang, Daniel M. Ziegler, Nisan Sti-	1281
1222	Reena Jana, Reza Ardeshtir Rokni, Rishabh Agar-	ennon, Ryan Lowe, Jan Leike, and Paul Christiano.	1282
1223	wal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy,	2021. Recursively summarizing books with human	1283
1224	Sarah Perrin, Sébastien M. R. Arnold, Sebastian	feedback. <i>Preprint</i> , arXiv:2109.10862.	1284
1225	Krause, Shengyang Dai, Shruti Garg, Shruti Sheth,		
1226	Sue Ronstrom, Susan Chan, Timothy Jordan, Ting	Sierra Wyllie, Ilia Shumailov, and Nicolas Papernot.	1285
1227	Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky,	2024. Fairness feedback loops: Training on synthetic	1286
1228	Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh	data amplifies bias. <i>Preprint</i> , arXiv:2403.07857.	1287
1229	Meshrum, Vishal Dharmadhikari, Warren Barkley,		
1230	Wei Wei, Wenming Ye, Woohyun Han, Woosuk	Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang,	1288
1231	Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan	Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang.	1289
1232	Wei, Victor Cotruta, Phoebe Kirk, Anand Rao,	2024. Iterative preference learning from human feed-	1290
1233	Minh Giang, Ludovic Peran, Tris Warkentin, Eli	back: Bridging theory and practice for rlhf under	1291
1234	Collins, Joelle Barral, Zoubin Ghahramani, Raia	kl-constraint. <i>Preprint</i> , arXiv:2312.11456.	1292
1235	Hadsell, D. Sculley, Jeanine Banks, Anca Dragan,		
1236	Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hass-	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,	1293
1237	abis, Koray Kavukcuoglu, Clement Farabet, Elena	Pu Zhao, Jiazhao Feng, Chongyang Tao, and Daxin	1294
1238	Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Ar-	Jiang. 2023. Wizardlm: Empowering large language	1295
1239	mand Joulin, Kathleen Kenealy, Robert Dadashi,	models to follow complex instructions. <i>Preprint</i> ,	1296
1240	and Alek Andreev. 2024. Gemma 2: Improving	arXiv:2304.12244.	1297
1241	open language models at a practical size. <i>Preprint</i> ,		
1242	arXiv:2408.00118.	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	1298
1243	Aman Singh Thakur, Kartik Choudhary, Venkat Srinik	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	1299
1244	Ramayapally, Sankaran Vaidyanathan, and Dieuwke	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	1300
1245	Hupkes. 2024. Judging the judges: Evaluating align-	ran Wei, Huan Lin, Jialong Tang, Jialin Wang,	1301
1246	ment and vulnerabilities in llms-as-judges. <i>Preprint</i> ,	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	1302
1247	arXiv:2406.12624.	Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai,	1303
		Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Ke-	1304
1248	Lewis Tunstall, Edward Beeching, Nathan Lambert,	qin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni,	1305
1249	Nazneen Rajani, Kashif Rasul, Younes Belkada,	Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize	1306
1250	Shengyi Huang, Leandro von Werra, Clémentine	Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan,	1307
1251	Fourrier, Nathan Habib, Nathan Sarrazin, Omar San-	Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge,	1308
1252	seviero, Alexander M. Rush, and Thomas Wolf. 2023.	Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren,	1309
1253	Zephyr: Direct distillation of lm alignment. <i>Preprint</i> ,	Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing	1310
1254	arXiv:2310.16944.	Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan,	1311
		Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang,	1312
1255	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa	Zhifang Guo, and Zhihao Fan. 2024. Qwen2 techni-	1313
1256	Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh	cal report. <i>Preprint</i> , arXiv:2407.10671.	1314
1257	Hajjishirzi. 2023. Self-instruct: Aligning language		
1258	models with self-generated instructions. In <i>Proceed-</i>	Michael JQ Zhang, Zhilin Wang, Jena D. Hwang,	1315
1259	<i>ings of the 61st Annual Meeting of the Association for</i>	Yi Dong, Olivier Delalleau, Yejin Choi, Eunsol Choi,	1316
1260	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Xiang Ren, and Valentina Pyatkin. 2024. Diverging	1317
1261	pages 13484–13508, Toronto, Canada. Association	preferences: When do annotators disagree and do	1318
1262	for Computational Linguistics.	models know? <i>Preprint</i> , arXiv:2410.14632.	1319
1263	Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-	Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie,	1320
1264	labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva	Yejin Choi, and Yuntian Deng. 2024. Wildchat:	1321
1265	Naik, Arjun Ashok, Arut Selvan Dhanasekaran,	lm chatgpt interaction logs in the wild. <i>Preprint</i> ,	1322
1266	Anjana Arunkumar, David Stap, Eshaan Pathak,	arXiv:2405.01470.	1323
1267	Giannis Karamanolakis, Haizhi Lai, Ishan Puro-		
1268	hit, Ishani Mondal, Jacob Anderson, Kirby Kuznia,	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	1324
1269	Krima Doshi, Kuntal Kumar Pal, Maitreya Patel,	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	1325
1270	Mehrad Moradshahi, Mihir Parmar, Mirali Purohit,	Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang,	1326
		Joseph E Gonzalez, and Ion Stoica. 2023a. Judging	1327

llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023b. *Judging llm-as-a-judge with mt-bench and chatbot arena*. Preprint, arXiv:2306.05685.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. *Llamafactory: Unified efficient fine-tuning of 100+ language models*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

A Prompts

A.1 Prompt for Feedback Signals Identification

The following is the full prompt we used for dialogue state tracking and SAT/DSAT classification. In addition, we also prompt GPT-4 to do domain and intent classification. The prompt is adapted from Das et al. (2023) and Lin et al. (2024b).

LABEL DEFINITION

```
{
  "valid_preceding_topical_relation_labels": [
    {
      "label": "YES",
      "definition": "The current turn has **some or any** topical/subtopical relation to the preceding conversation context."
    },
    {
      "label": "NO",
      "definition": "The current turn has **absolutely no** topical/subtopical relation to the preceding conversation context OR is the first turn in the conversation, marking the beginning of a new dialogue segment."
    }
  ],
  "valid_domain_labels": [
    "AI MACHINE LEARNING AND DATA SCIENCE",
    "ASTROLOGY",
    "BIOLOGY AND LIFE SCIENCE",
    "BUSINESS AND MARKETING",
```

```
"CAREER AND JOB APPLICATION",
"CLOTHING AND FASHION",
"COOKING FOOD AND DRINKS",
"CRAFTS",
"CULTURE AND HISTORY",
"CYBERSECURITY",
"DATING FRIENDSHIPS AND RELATIONSHIPS",
"DESIGN",
"EDUCATION",
"ENTERTAINMENT",
"ENVIRONMENT AGRICULTURE AND ENERGY",
"FAMILY PARENTING AND WEDDINGS",
"FINANCE AND ECONOMICS",
"GAMES",
"GEOGRAPHY AND GEOLOGY",
"HEALTH AND MEDICINE",
"HOUSING AND HOMES",
"HUMOR AND SARCASM",
"LANGUAGE",
"LAW AND POLITICS",
"LITERATURE AND POETRY",
"MANUFACTURING AND MATERIALS",
"MATH LOGIC AND STATISTICS",
"MUSIC AND AUDIO",
"NEWS",
"PETS AND ANIMALS",
"PHILOSOPHY",
"PHYSICS CHEMISTRY AND ASTRONOMY",
"PRODUCTIVITY",
"PSYCHOLOGY AND EMOTIONS",
"RELIGION AND MYTHOLOGY",
"SHIPPING AND DELIVERY",
"SHOPPING AND GIFTS",
"SMALL TALK",
"SOCIAL MEDIA",
"SOFTWARE AND WEB DEVELOPMENT",
"SPORTS AND FITNESS",
"TAXATION",
"TECHNOLOGY",
"TIME AND DATES",
"TRANSPORTATION AUTOMOTIVE AND AEROSPACE",
"TRAVEL",
"VISUAL ARTS AND PHOTOGRAPHY",
"WEATHER",
"WRITING JOURNALISM AND PUBLISHING",
"OTHER"
],
  "valid_intent_labels": [
    {
      "label": "INTENT:1-INFORMATION-SEEKING",
      "definition": "The user wants to find
```

1432	factual information or answers to	applicable"	1484
1433	specific questions."	},	1485
1434	},	{	1486
1435	{	"label": "Praise",	1487
1436	"label": "INTENT:2-ANALYSIS",	"definition": "The user uses positive	1488
1437	"definition": "The user asks analytical	feedback words (e.g., excellent, amazing)	1489
1438	or conceptual questions about a	or emojis, indicating enthusiasm and	1490
1439	complex topic or problem. The user's	enjoyment of the conversation"	1491
1440	questions require some degree of	},	1492
1441	reasoning, interpretation, argumentation,	{	1493
1442	comparison, and/or data processing."	"label": "Personal_Details",	1494
1443	},	"definition": "The user shares more	1495
1444	{	personal details or opinions with the AI	1496
1445	"label": "INTENT:3-CREATION",	agent when satisfied with its responses"	1497
1446	"definition": "The user asks the agent	},	1498
1447	to either generate original content	{	1499
1448	or translate existing content into new	"label": "Humor",	1500
1449	content based on specified criteria or	"definition": "The user jokes with or	1501
1450	constraints."	challenges the AI agent in a friendly	1502
1451	},	manner when suitable"	1503
1452	{	},	1504
1453	"label": "INTENT:4-OPEN-ENDED_DISCOVERY",	{	1505
1454	"definition": "The user wants to	"label": "Acknowledgment",	1506
1455	casually chat or play with the	"definition": "The user acknowledges or	1507
1456	agent out of curiosity, boredom,	confirms that they understood or agreed	1508
1457	or humor, OR the user's intent is	with the AI agent's explanations when	1509
1458	so unclear/underspecified that it's	relevant"	1510
1459	impossible to categorize in any of the	},	1511
1460	other intent classes. The user mainly	{	1512
1461	treats the agent as a conversation or	"label": "Positive_Closure",	1513
1462	chitchat partner, and none of the other	"definition": "The user ends the	1514
1463	intent categories can be assigned."	conversation on a positive note	1515
1464	}	without asking for more information	1516
1465	},	or assistance"	1517
1466	"valid_satisfaction_labels": [},	1518
1467	{	{	1519
1468	"label": "Gratitude",	"label": "Getting_There",	1520
1469	"definition": "The user thanks or	"definition": "The user acknowledges that	1521
1470	compliments the AI agent for its	the model's response is getting better	1522
1471	responses"	or has merit but is not fully satisfied.	1523
1472	},	Appropriate dissatisfaction criteria	1524
1473	{	may need to be checked as well when	1525
1474	"label": "Learning",	Getting_There presents"	1526
1475	"definition": "The user learns something	},	1527
1476	new or useful by indicating curiosity	{	1528
1477	and satisfaction with the information	"label": "N/A",	1529
1478	provided"	"definition": "The user utterance of the	1530
1479	},	turn does NOT match the definition of	1531
1480	{	any other valid satisfaction labels"	1532
1481	"label": "Compliance",	}	1533
1482	"definition": "The user follows the AI	},	1534
1483	agent's suggestions or instructions when	"valid_dissatisfaction_labels": [1535

1536	{	provided by the AI agent"	1588
1537	"label": "Negative_Feedback",	},	1589
1538	"definition": "The user explicitly	{	1590
1539	expresses dissatisfaction, frustration,	"label": "Style",	1591
1540	annoyance, or anger with the AI agent's	"definition": "The user feels that there	1592
1541	response or behavior"	is a mismatch between their preferred	1593
1542	},	style (e.g. bullet point vs paragraph,	1594
1543	{	formal vs casual, short vs long, etc.)	1595
1544	"label": "Revision",	and what is provided by the AI agent"	1596
1545	"definition": "The user explicitly asks	},	1597
1546	the AI agent to revise its previous	{	1598
1547	response or repeatedly asks similar	"label": "N/A",	1599
1548	questions"	"definition": "The user utterance of the	1600
1549	},	turn does NOT match the definition of	1601
1550	{	any other valid dissatisfaction labels"	1602
1551	"label": "Factual_Error",	}	1603
1552	"definition": "The user points out the AI],	1604
1553	agent's factual mistakes, inaccuracies,	"valid_state_labels": [1605
1554	or self-contradiction in its information	{	1606
1555	or output"	"label": "FEEDBACK",	1607
1556	},	"definition": "The user utterance of the	1608
1557	{	turn contains a comment or evaluation or	1609
1558	"label": "Unrealistic_Expectation",	judgement of the previous turn's agent	1610
1559	"definition": "The user has unrealistic	response"	1611
1560	expectations of what the AI agent can do	},	1612
1561	and does not accept its limitations or	{	1613
1562	alternatives"	"label": "REFINEMENT",	1614
1563	},	"definition": "The user utterance of the	1615
1564	{	turn is a repetition or refinement of	1616
1565	"label": "No_Engagement",	unclear/underspecified instruction given	1617
1566	"definition": "The user does not respond	in the previous turn's user utterance"	1618
1567	to the AI agent's questions, suggestions,	},	1619
1568	feedback requests, etc."	{	1620
1569	},	"label": "NEWTOPIC",	1621
1570	{	"definition": "The user utterance of the	1622
1571	"label": "Ignored",	turn is either the first turn of the	1623
1572	"definition": "The user implies that	conversation or is not related in terms	1624
1573	their query was ignored completely or	of topic or task to its previous turn,	1625
1574	that the response did not address their	introducing a new topic or task"	1626
1575	intent/goal at all"	},	1627
1576	},	{	1628
1577	{	"label": "CONTINUATION",	1629
1578	"label": "Lower_Quality",	"definition": "The user utterance of the	1630
1579	"definition": "The user perceives a	turn is a topical or logical continuation	1631
1580	decline in quality of service compared	of the previous turn"	1632
1581	to previous experience with other	}	1633
1582	agents/tools, etc."]	1634
1583	},	}	1635
1584	{		1636
1585	"label": "Insufficient_Detail",	## TASK ##	1637
1586	"definition": "The user wants more	You are given a dialogue between a user and an	1638
1587	specific/useful information than what is	agent comprised of turns starting with T. For each	1639

turn, solely based on the turn’s User utterance, you must carefully analyze the conversation and answer the following questions by replacing \$instruction\$ with correct answers in JSON format. - Summarize the user utterance in ≤ 3 sentences

- Analyze the user utterance’s relation with the previous turn and output an appropriate label from the “valid_preceding_topical_relation_labels” list.
- Analyze the user utterance’s domain and output an appropriate label from the “valid_domain_labels” list. If preceding_topical_relation is YES, the domain label must be consistent with the preceding turn’s domain label.
- Analyze the user utterance’s intent and output an appropriate label from the “valid_intent_labels” list.
- Analyze the user utterance’s satisfaction with respect to the previous turn’s AI response and output all applicable labels from the “valid_satisfaction_labels” list.
- Analyze the user utterance’s dissatisfaction with respect to the previous turn’s AI response and output all applicable labels from the “valid_dissatisfaction_labels” list.
- Analyze the user utterance’s state and output an appropriate label from the “valid_state_labels” list.

OUTPUT FORMAT

The length and turn order of the output list must match the length and turn order of the input list. The sample output format is given as follow: [{

```
"T-$turn number$": {
  "summary": "$turn summary in  $\leq 3$  sentence$",
  "preceding_topical_relation": "$an appropriate valid preceding topical relation label$",
  "domain": "$an appropriate valid domain label$",
  "intent": "INTENT:$an appropriate valid intent label$",
  "satisfaction": [$a comma separated string list of applicable valid satisfaction label(s)$],
  "dissatisfaction": [$a comma separated string list of applicable valid dissatisfaction label(s)$],
  "state": "$an appropriate valid state label$"
}
```

}]

INPUT

#D1#

OUTPUT

A.2 Prompt for Preference Pair Construction

The following is the prompt for constructing preference data.

Conversation between User and AI

< |begin_of_history| >

history

< |end_of_history| >

Instruction

What are the user’s query and preferences? The query should be the user’s first attempt before providing any feedbacks to the model. Only output the turn id. The preference should always be based on user’s feedbacks and in complete sentences. Generate your answer in json format like

```
[ {
  "query": turn id,
  "preferences": [preference 1, preference 2, ...]
} ]
```

A.3 Prompt for Checklist-guided Evaluation

The following is the prompt for checklist-guided evaluation. We borrow the WB-Reward prompt from WILDBENCH (Lin et al., 2024a).

Instruction

You are an expert evaluator. Your task is to evaluate the quality of the responses generated by two AI models. We will provide you with the user query and a pair of AI-generated responses (Response A and B). You should first read the user query and the conversation history carefully for analyzing the task, and then evaluate the quality of the responses based on and rules provided below.

Conversation between User and AI

History

< |begin_of_history| >

{history}

< |end_of_history| >

Current User Query

< |begin_of_query| >

{query}

< |end_of_query| >

Response A

< |begin_of_response_A| >

{response_a}

< |end_of_response_A| >

```

1742 ## Response B
1743 < |begin_of_response_B| >
1744 {response_b}
1745 < |end_of_response_B| >
1746 # Evaluation
1747 ## Checklist
1748 < |begin_of_checklist| >
1749 {checklist}
1750 < |end_of_checklist| >
1751 Please use this checklist to guide your evaluation,
1752 but do not limit your assessment to the checklist.
1753 ## Rules
1754 You should compare the above two responses
1755 based on your analysis of the user queries and the
1756 conversation history. You should first write down
1757 your analysis and the checklist that you used for
1758 the evaluation, and then provide your assessment
1759 according to the checklist. There are five choices
1760 to give your final assessment: ["A++", "A+",
1761 "A=B", "B+", "B++"], which correspond to the
1762 following meanings:
1763 - 'A++': Response A is much better than Response
1764 B.
1765 - 'A+': Response A is only slightly better than
1766 Response B.
1767 - 'A=B': Response A and B are of the same quality.
1768 Please use this choice sparingly.
1769 - 'B+': Response B is only slightly better than
1770 Response A.
1771 - 'B++': Response B is much better than Response
1772 A.
1773 ## Output Format
1774 First, please output your analysis for each model
1775 response, and then summarize your assessment
1776 to three aspects: "reason A=B", "reason A >
1777 B", and "reason B > A", and finally make your
1778 choice for the final assessment. Please provide
1779 your evaluation results in the following json format
1780 by filling in the placeholders in []:
1781 {
1782 "analysis of A": "[analysis of Response
1783 A]",
1784 "analysis of B": "[analysis of Response
1785 B]",
1786 "reason of A=B": "[where Response A and
1787 B perform equally well]",
1788 "reason of A>B": "[where Response A is
1789 better than Response B]",
1790 "reason of B>A": "[where Response B is
1791 better than Response A]",
1792 "choice": "[A++ or A+ or A=B or B+ or
1793 B++]"
```

```

}
```

A.4 Prompt for Dataset Evaluation

The following is the prompt for constructing the on-policy version of the ULTRAFeEDBACK dataset. The prompt is adapted from the WB-Reward prompt (Lin et al., 2024a).

Instruction

You are an expert evaluator. Your task is to evaluate the quality of the responses generated by two AI models. We will provide you with the user query and a set of AI-generated responses (Response A, Response B, Response C, Response D, Response E). You should first read the user query and the conversation history carefully for analyzing the task, and then evaluate the quality of the responses based on the rules provided below.

Conversation between User and AI

History

```
< |begin_of_history| >
```

```
{history}
```

```
< |end_of_history| >
```

Current User Query

```
< |begin_of_query| >
```

```
{query}
```

```
< |end_of_query| >
```

Response A

```
< |begin_of_response_A| >
```

```
{response_a}
```

```
< |end_of_response_A| >
```

Response B

```
< |begin_of_response_B| >
```

```
{response_b}
```

```
< |end_of_response_B| >
```

Response C

```
< |begin_of_response_C| >
```

```
{response_c}
```

```
< |end_of_response_C| >
```

Response D

```
< |begin_of_response_D| >
```

```
{response_d}
```

```
< |end_of_response_D| >
```

Response E

```
< |begin_of_response_E| >
```

```
{response_e}
```

```
< |end_of_response_E| >
```

Evaluation

Checklist

```
< |begin_of_checklist| >
```

```
{checklist}
```

```
< |end_of_checklist| >
```

Please use this checklist to guide your evaluation, but do not limit your assessment to the checklist.

Rules

You should compare the above five responses based on your analysis of the user queries and the conversation history. You should first write down your analysis and the checklist that you used for the evaluation, and then provide your assessment according to the checklist.

There are six choices to give your final assessment: ["A", "B", "C", "D", "E", "A=B=C=D=E"], which correspond to the following meanings:

- 'A': Response A is much better than the other responses.
- 'B': Response B is much better than the other responses.
- 'C': Response C is much better than the other responses.
- 'D': Response D is much better than the other responses.
- 'E': Response E is much better than the other responses.
- 'A=B=C=D=E': Response A, B, C, D, E are of the same quality. No response particularly stood out. Please use this choice sparingly.

Output Format

First, please output your analysis for each model response, and then summarize your assessment to "comparison of A, B, C, D, E", and finally make your choice for the final assessment. Please provide your evaluation results in the following json format by filling in the placeholders in []:

```
{
  "analysis of A": "[analysis of Response A]",
  "analysis of B": "[analysis of Response B]",
  "analysis of C": "[analysis of Response C]",
  "analysis of D": "[analysis of Response D]",
  "analysis of E": "[analysis of Response E]",
  "comparison of A, B, C, D, E": "[where Response A, B, C, D, E perform equally well]",
  "choice": "[A or B or C or D or E or A=B=C=D=E]"
}
```

B SAT and DSAT

B.1 Detailed SAT and DSAT Criteria

The detailed definitions of SAT and DSAT can be found in Table 4 and Table 5.

B.2 SAT and DSAT Annotation

We randomly sampled 50 multi-turn conversations, totaling over 500 utterances, and assigned 4 expert annotators to perform the same classification task. Each conversation was annotated by at least 2 annotators, resulting in a final Cohen's Kappa agreement of $\kappa = 0.70$ for SAT and $\kappa = 0.54$ for DSAT. For human annotation, we utilized a web-based annotation tool named Potato (Pei et al., 2022). The interface is shown in Figure 4. After completing the annotations, the annotators reviewed and discussed any disagreements, resolving conflicts to establish a ground truth test set of 50 conversations. GPT-4's performances on SAT and DSAT classification can be found in table 6. GPT-4 demonstrates strong performance in classifying SAT (satisfaction) signals, with high accuracy at 91.7% and balanced precision and recall, both around 73%. The Cohen's Kappa of 68.5% reflects substantial agreement with human annotators. For DSAT (dissatisfaction) signals, GPT-4 achieves a precision of 83.3%, with a recall of 48.4%, leading to an F1 score of 61.2% and a Cohen's Kappa of 50.4%. These metrics indicate that GPT-4 is effective at recognizing both SAT and DSAT signals.

C GPT-4's Performance on Checklist-guided Evaluation

We randomly selected 200 multi-turn conversations, and assigned 6 expert annotators to perform checklist-guided evaluation. Each conversation is annotated by at least 2 annotators, resulting in a final Cohen's Kappa agreement of $\kappa = 43.6$. After completing the annotations, the annotators reviewed and discussed any disagreements, resolving conflicts to establish a ground truth test set. For human annotation, we utilized a web-based annotation tool named Potato (Pei et al., 2022). The interface is shown in Figure 5. GPT-4's performances on checklist-guided evaluation can be found in Table 7. Our findings indicate that GPT-4's ability to perform checklist-guided evaluation has a relatively high agreement with human annotators, achieving a Cohen's Kappa of $\kappa = 37.2$. GPT-4 performs relatively on par with humans on checklist-guided evaluation.

Keyword	Definition
Gratitude	The user thanks or compliments the AI agent for its responses.
Learning	The user learns something new or useful by indicating curiosity and satisfaction with the information provided.
Compliance	The user follows the AI agent’s suggestions or instructions when applicable.
Praise	The user uses positive feedback words (e.g., excellent, amazing) or emojis, indicating enthusiasm and enjoyment of the conversation.
Personal Details	The user shares more personal details or opinions with the AI agent when satisfied with its responses.
Humor	The user jokes with or challenges the AI agent in a friendly manner when suitable.
Acknowledgment	The user acknowledges or confirms that they understood or agreed with the AI agent’s explanations when relevant.
Positive Closure	The user ends the conversation on a positive note without asking for more information or assistance.
Getting There	The user acknowledges that the model’s response is getting better or has merit but is not fully satisfied.

Table 4: Detailed definitions of the SAT Rubrics.

D Implementation Details

We found that hyperparameter tuning is crucial for achieving optimal performance in preference optimization. Generally, on-policy data requires a lower learning rate than GPT-4o data, and instruct models need a lower learning rate than base models. Specifically, Mistral and Gemma (Team et al., 2024) require a lower learning rate than Phi 3, LLaMA 3 and Qwen 2. Initially, we followed the Zephyr setup (Tunstall et al., 2023), which employs a learning rate of $2e-5$ for supervised fine-tuning (SFT). However, we found that our models quickly collapsed, failing to generate sensible outputs after just a few dozen iterations. After conducting a grid search on the hyperparameters for both SFT and DPO training, we discovered that while it is acceptable to use a larger learning rate for training base models, a much smaller learning rate is required for instruct models, likely due to the various annealing techniques applied during the post-training process (Parmar et al., 2024). We also explored NLL regularization (Liu et al., 2024a) with a regularization strength of 0.2, but the results are not ideal, and therefore, we did not include NLL regularization in the final set up. We trained all the models using LLaMA Factory (Zheng et al., 2024), a unified efficient LLM finetuning framework. LLaMA Factory

is licensed under the Apache-2.0 License. The following is the hyperparameters we used in our final experiment.

SFT Training. For SFT training, we trained all the models for 1 epoch with a batch size of 128, a learning rate of $5e-6$, a linear warm-up ratio of 0.1, and a cosine learning rate scheduler. Additionally, it is recommended to use a higher learning rate (e.g., $2e-5$) if you are fine-tuning from the base models. It takes about 8 A100 GPU hours to finish.

DPO Training. For DPO training, we trained all the models for 1 epoch with a batch size of 32, a learning rate of $5e-7$, and $\beta = 0.1$. All other hyperparameters remained the same as in the SFT training. It takes about 24 A100 GPU hours to finish.

E WildChat Dataset

The WildChat Dataset is a corpus of 1 million real-world user-ChatGPT interactions, covering a wide range of languages and user prompts. Most of the conversations are single-turn. It was constructed by offering free access to ChatGPT and GPT-4 in exchange for consensual chat history collection and is licensed under the Open Data Commons Attribution License (ODC-By) v1.0. To protect

Keyword	Definition
Negative Feedback	The user explicitly expresses dissatisfaction, frustration, annoyance, or anger with the AI agent’s response or behavior.
Revision	The user explicitly asks the AI agent to revise its previous response or repeatedly asks similar questions.
Factual Error	The user points out the AI agent’s factual mistakes, inaccuracies, or self-contradiction in its information or output.
Unrealistic Expectation	The user has unrealistic expectations of what the AI agent can do and does not accept its limitations or alternatives.
No Engagement	The user does not respond to the AI agent’s questions, suggestions, feedback requests, etc.
Ignored	The user implies that their query was ignored completely or that the response did not address their intent/goal at all.
Lower Quality	The user perceives a decline in quality of service compared to previous experience with other agents/tools, etc.
Insufficient Detail	The user wants more specific/useful information than what is provided by the AI agent.
Style	The user feels that there is a mismatch between their preferred style and what is provided by the AI agent.

Table 5: Detailed definitions of the DSAT Rubrics.

	Accuracy	Precision	Recall	F1	GPT-Human κ	Human-Human κ
SAT	91.7	73.2	73.6	73.4	68.5	70.0
DSAT	81.8	83.3	48.4	61.2	50.4	54.1

Table 6: Agreement on SAT and DSAT Classification. All numbers are in %.

personally identifiable information (PII), WildChat employed Microsoft’s Presidio⁴ as the framework, SpaCy⁵ for Named Entity Recognition, and custom rules to remove PII—including names, phone numbers, emails, credit cards, and URLs—across multiple languages such as English, Chinese, Russian, French, Spanish, German, Portuguese, Italian, Japanese, and Korean. Additionally, WildChat utilized GeoLite2⁶ to map IP addresses to countries and states before hashing them for privacy. While WildChat releases only hashed IP addresses and request headers (including browser details and accepted languages), these identifiers could allow researchers to infer connections between conversations from the same user, though no direct linkage is provided in the dataset.

⁴<https://microsoft.github.io/presidio/>

⁵<https://spacy.io/>

⁶<https://dev.maxmind.com/geoip/geo-lite2-free-geolocation-data>

Conversation-32

Turn-1
User: What is a semaphore?
Agent: As an AI language model, I don't know the term "sephamore." It may be a misspelling of "semaphore," which is a system of sending signals by using flags, lights, or mechanical devices to communicate messages over long distances. Semaphores are commonly used in naval or military operations, railway signaling systems, and computer programming for synchronizing multiple processes.

Turn-2
User: Yes I meant semaphore
Agent: Semaphore is a system of communication that uses a combination of flags or lights to send messages over relatively long distances. The flags or lights are usually arranged in a specific way to convey a specific message, allowing multiple messages to be conveyed quickly and efficiently. Semaphore is commonly used in military (naval) and railway signaling systems, air traffic control signals, and computer programming to manage concurrent processes.

Turn-3
User: Explain it in the context of programming
Agent: In programming, a semaphore is a synchronization mechanism that is used to manage access to a shared resource by multiple processes or threads in a concurrent or parallel processing system. It acts as a lock on a resource, which ensures that only one process or thread can access the resource at any given time. A semaphore contains a counter that is initially set to a specific value that represents the maximum number of processes or threads that can simultaneously access the shared resource. When a process or thread requests access to the resource, the semaphore decrements the counter, and if the counter reaches zero, the semaphore blocks the request and the process or thread waits until the resource is released by another process or thread. Once a process or thread has finished accessing the resource, it releases it and increments the semaphore counter, which allows another waiting process or thread to access the resource. This way, semaphores can be used to ensure mutual exclusion and synchronization between multiple processes or threads that need to access shared resources in a concurrent processing system.

Turn-4
User: Write one in C

Select an appropriate State

☐ New Topic
☒ Feedback
☐ Refinement
☐ Continuation

Select all applicable Satisfaction criteria

☒ N/A

☐ 1. **Gratitude:** The user thanks or compliments the AI agent for its responses
☐ 2. **Learning:** The user learns something new or useful by indicating curiosity and satisfaction with the information provided
☐ 3. **Compliance:** The user follows the AI agent's suggestions or instructions when applicable
☐ 4. **Praise:** The user uses positive feedback words (e.g., excellent, amazing) or emojis, indicating enthusiasm and enjoyment of the conversation
☐ 5. **Personal Details:** The user shares more personal details or opinions with the AI agent when satisfied with its responses
☐ 6. **Humor:** The user jokes with or challenges the AI agent in a friendly manner when suitable
☐ 7. **Acknowledgment:** The user acknowledges or confirms that they understood or agreed with the AI agent's explanations when relevant
☐ 8. **Positive Closure:** The user ends the conversation on a positive note without asking for more information or assistance
☐ 9. **Getting There:** The user acknowledges that the model's response is getting better or has merit but is not fully satisfied

Select all applicable Dissatisfaction criteria

☐ N/A

☐ 1. **Negative Feedbacks:** The user explicitly expresses dissatisfaction, frustration, annoyance, or anger with the AI agent's response or behavior
☐ 2. **Revision:** The user explicitly asks the AI agent to revise its previous response or repeatedly asks similar questions
☐ 3. **Factual Errors:** The user points out the AI agent's factual mistakes, inaccuracies, or self-contradiction in its information or output
☐ 4. **Unrealistic Expectation:** The user has unrealistic expectations of what the AI agent can do and does not accept its limitations or alternatives
☐ 5. **No Engagement:** The user does not respond to the AI agent's questions, suggestions, feedback requests, etc
☐ 6. **Ignored:** The user implies that their query was ignored completely or that the response did not address their intent/goal at all
☐ 7. **Lower Quality:** The user perceives a decline in quality of service compared to previous experience with other agents/tools, etc
☒ 8. **Insufficient Detail:** The user wants more specific/useful information than what is provided by the AI agent
☐ 9. **Style:** The user feels that there is a mismatch between their preferred style (long vs short, bullet point vs paragraph, etc.) and what is provided by the AI agent

Figure 4: The interface used for annotating SAT and DSAT signals.

GPT-Human κ	Human-Human κ	GPT-Human Agreement	Human-Human Agreement
37.2	43.6	57.14	63.27

Table 7: Agreement on checklist-guided Evaluation. All numbers are in %.

Dialogue Analysis
Home
Statistics
Help
Finished 47/50
Current_id 11

Currently logged in as timcook@apple.com

Turn-1

User:

Is MrBeast alive?

Checklist:

The response should be safe.

The user prefers concise and direct answers without additional information.

Response A:

Yes, MrBeast is alive.

Response B:

As of my last update in October 2023, MrBeast, whose real name is Jimmy Donaldson, is alive. However, for the most current information, I recommend checking the latest news or his social media profiles.

Judge the quality of the model responses.

☒ A++: Response A is much better than Response B.
☐ A+: Response A is only slightly better than Response B.
☐ A=B: Response A and B are of the same quality. Please use this choice sparingly.
☐ B+: Response B is only slightly better than Response A.
☐ B++: Response B is much better than Response A.

Keyboard Input:

←

→

1

Quality: A++: Response A is much better than Response B.

2

Quality: A+: Response A is only slightly better than Response B.

3

Quality: A=B: Response A and B are of the same quality. Please use this choice sparingly.

4

Quality: B+: Response B is only slightly better than Response A.

5

Quality: B++: Response B is much better than Response A.

Move backward

Move forward

Figure 5: The interface used for annotating checklist-guided evaluation.