

# THE ERA OF REAL-WORLD HUMAN INTERACTION: RL FROM USER CONVERSATIONS

Anonymous authors

Paper under double-blind review

## ABSTRACT

We posit that to achieve continual model improvement and multifaceted alignment, future models must learn from natural human interaction. Current conversational models are aligned using pre-annotated, expert-generated human feedback. In this work, we introduce Reinforcement Learning from Human Interaction (RLHI), a post-training paradigm that learns directly from in-the-wild user conversations. We develop two complementary methods: (1) *RLHI with User-Guided Rewrites*, which revises unsatisfactory model outputs based on users’ natural-language follow-up responses, (2) *RLHI with User-Based Rewards*, which learns via a reward model conditioned on knowledge of the user’s long-term interaction history (termed persona). Together, these methods link long-term user personas to turn-level preferences via persona-conditioned preference optimization. Trained on conversations derived from WildChat, both RLHI variants outperform strong baselines in personalization and instruction-following, and similar feedback enhances performance on reasoning benchmarks. These results suggest organic human interaction offers scalable, effective supervision for personalized alignment.

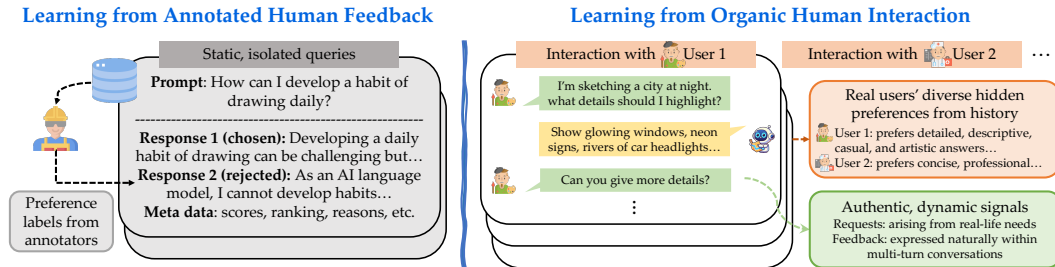


Figure 1: **From annotated feedback to the era of real-world human interaction.** **Left:** Traditional alignment relies on expert-curated annotations of ranked responses or labels, providing static, out-of-distribution supervision. **Right:** In-the-wild conversations reveal users’ long-term histories, dynamic demands, and diverse signals, enabling personalized, contextual, and continual learning.

## 1 INTRODUCTION

Today, language model post-training primarily depends on static corpora of expert-annotated data: verifiable questions, fixed demonstrations, and rankings or ratings collected outside of natural conversational contexts. While these datasets are effective for instilling general capabilities, they reflect the opinions and heuristics of annotators in unnatural scenarios rather than the *authentic, diverse long-term goals and preferences of real users*; they capture static, context-free judgments instead of *evolving, situational demands*; and they scale with labeling budgets rather than with actual usage and diversity of organic users, as is illustrated on the left side of Figure 1.

In contrast, humans learn and improve through continual experience by interacting with their environment and other actors, receiving feedback, and adjusting behavior over time (Tomasello et al., 2005). Likewise, a rich and organic source of supervision for language models already exists in the wild: **human interaction**—the ongoing, natural exchanges between models and real users. As is

---

054 shown on the right side of Figure 1, such organic interactions reveal hidden user preferences from  
055 long-term histories and dynamic, context-dependent demands, as people reveal their priorities and  
056 concerns not through annotation formats, but by discussing what matters to them, revising or re-  
057 attempting questions, explicitly or implicitly approving or critiquing model outputs, following up,  
058 or switching goals mid-dialogue. Because they arise directly from model outputs in authentic usage  
059 contexts, such interactions provide a rich signal for learning personalized and adaptable behavior,  
060 paving the way toward personal superintelligence. While this source of supervision has historically  
061 been hard to extract, resulting in resorting to collecting static training data instead, the power of  
062 modern language models now gives us a greater ability to extract these signals.

063 To achieve this vision, we introduce RLHI, a post-training paradigm that learns directly from in-  
064 the-wild conversations through two complementary methods: (1) *RLHI with User-Guided Rewrites*  
065 (§2.3), which revises unsatisfactory model outputs based on users’ natural-language follow-ups, and  
066 pairs the rewrites with the originals for preference learning; and (2) *RLHI with User-Based Rewards*  
067 (§2.4), which ranks candidate responses using a reward model conditioned on user personas derived  
068 from long-term histories to generate preference pairs. Together, these methods link long-term user  
069 personas to turn-level preferences via persona-conditioned preference optimization.

070 We evaluate RLHI in three settings. (i) *User-based evaluation* with our WILDCHAT USEREVAL:  
071 both RLHI variants outperform strong baselines in personalization and instruction-following, and  
072 a human study corroborates these trends. (ii) *Standard instruction-following benchmarks*: *User-*  
073 *Based Rewards* attains a 77.9% length-controlled win rate on AlpacaEval 2.0, surpassing all RLHF  
074 methods. (iii) *Reasoning*: *User-Guided Rewrites* raises average accuracy from 26.5 to 31.8 across  
075 four benchmarks. Our ablation studies further show that RLHI benefits from user guidance and  
076 interaction diversity, that reinforcement learning outperforms supervised finetuning, and that quality  
077 filtering is essential for effectively leveraging noisy human interaction data.

## 078 079 080 2 RLHI: REINFORCEMENT LEARNING FROM HUMAN INTERACTION

### 081 082 2.1 THE ERA OF REAL-WORLD HUMAN INTERACTION

083  
084 Artificial intelligence (AI) has progressed rapidly in recent years through large-scale pretraining and  
085 fine-tuning with human examples and preferences. Yet this trajectory is slowing: high-quality data  
086 is running out, and imitation alone cannot push systems beyond existing human knowledge. Recent  
087 proposals call for an *era of experience* (Silver & Sutton, 2025), in which AI systems advance by  
088 continually learning from their own interactions with the world. Since these systems ultimately  
089 exist to assist humans, interaction with users becomes a natural and essential dimension of this shift.  
090 The *era of real-world human interaction* thus forms a core pillar of the era of experience, providing  
091 both the raw data and personalization signals necessary for adaptive, human-centered intelligence.

092 We define learning from human interaction as the process of improving AI models through natural,  
093 continual exchanges with real users. Such interactions may involve messages, actions, requests,  
094 or demonstrations provided in direct response to the model’s outputs. These exchanges not only  
095 reveal user goals and preferences but also create an evolving feedback loop that enables systems  
096 to refine their behavior over time. To truly benefit from human interaction, AI needs to go beyond  
097 coarse binary labels to absorb knowledge, preferences, reasoning skills, perceptual cues, cooperative  
098 strategies, and social norms, learning deeper forms of intelligence through interaction.

099 Compared with other training data sources, human interaction is distinguished by three key proper-  
100 ties: (1) **Contextual grounding** — arises within the flow of ongoing tasks or conversations, directly  
101 tied to the user’s situational needs and the model’s prior outputs, while being shaped by personal-  
102 ized knowledge of the user’s profile, history, and preferences; (2) **Evolving distribution** — reflects  
103 goals that shift, environments that change, and preferences that adapt over time, thereby providing  
104 supervision that is temporally relevant and aligned with the real distribution of human needs and  
105 priorities; and (3) **Diverse supervision signals** — appears in both explicit high-bandwidth signals  
106 beyond scalar rewards (e.g., corrections or clarifications) and implicit cues (e.g., disengagement or  
107 frustration), and may include style and role assignments, emotional tone, or even adversarial inputs  
such as jailbreak attempts, which require careful handling, but also offer valuable information.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

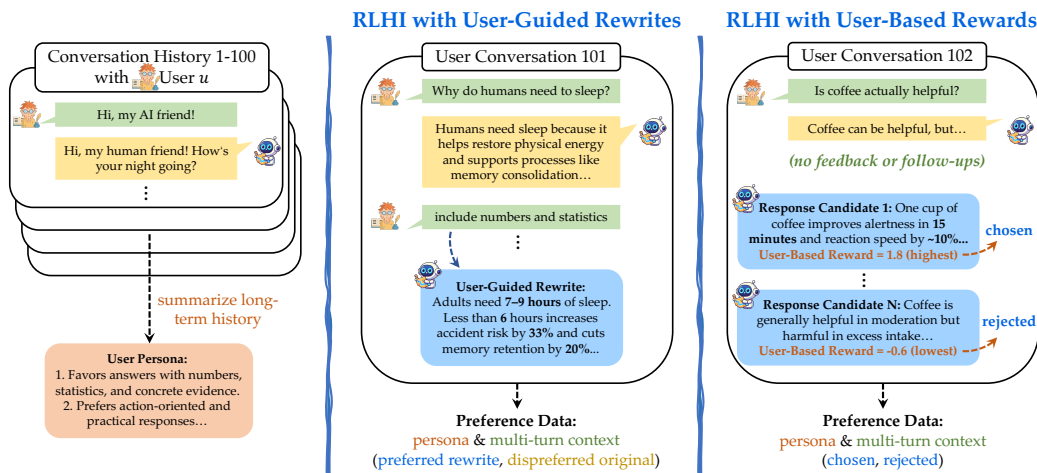


Figure 2: **Reinforcement Learning from Human Interaction (RLHI)**. We derive a natural-language persona summary from each user’s long-term conversational history. For real-world requests, RLHI operates in two modes: (1) User-Guided Rewrites, where unsatisfactory model outputs are revised based on users’ natural-language follow-ups, creating preference pairs between the original and rewritten responses; and (2) User-Based Rewards, where multiple candidate responses are generated and ranked by a reward model conditioned on the user’s persona, yielding chosen-rejected pairs. Both methods leverage personas and multi-turn context to enable personalized alignment.

In this paper, we focus on large language models that engage daily with millions of users. Here, human interaction takes the minimal form of textual messages, yet still conveys contextual, dynamic, and diverse requests, holding unique potential as a driver of continual model improvement.

## 2.2 ANALYSIS OF REAL-WORLD HUMAN INTERACTION

To determine the feasibility of our approach, we first consider *currently available human interaction data*, analyzing its properties. We note that these properties are necessarily tied to the capabilities of current models, and we expect these statistics to change considerably in the coming years.

**Users often provide feedback to improve model responses.** We analyze user messages in the WildChat-1M dataset, which contains over one million conversations with ChatGPT (Zhao et al., 2024b). In each multi-turn conversation, the first message is the *initial request*, and we prompt an GPT-4o model to classify user follow-up messages into four types: (1) *new requests*, where the user shifts to a new topic, substantially reformulates the original, or provides unrelated input; (2) *re-attempts with feedback*, where the user refines the initial prompt, adds clarification, or provides explicit or implicit feedback; (3) *re-attempts without feedback*, where the same prompt is repeated with no new input; and (4) *positive feedback*, where the user expresses praise or satisfaction. We find the distributions are: 27.07% of user messages are *initial requests*, 40.40% are *new requests*, 26.51% are *re-attempts with feedback*, 4.77% are *re-attempts without feedback*, and 1.25% are *positive feedback*, with more details and examples in Appendix A. Conversations of later stages are dominated by *re-attempts with feedback*, accounting for 83.15% of user utterances after the fifth turn. *re-attempts with feedback* are relatively short, averaging 272 characters compared to 725 for initial requests, but are semantically dense. We note that given the huge amount of human interactions in current production systems, these percentages convert to very large amounts of supervisory data. We note that while these are current statistics, in the future, as models display further capabilities, users will change their behavior. For example, if users know that models will learn from their textual feedback, then they are even more likely to provide it.

**Real-world human interaction data are more diverse than existing preference datasets.** Conversation messages span a wide range of forms and topics (for example, creative writing, analysis, and coding) and occur in conversations of highly varying length (average 2.54 turns). To quantify this diversity, we compare request contexts in our generated preference dataset with two widely

used annotated feedback datasets: HH-RLHF (Bai et al., 2022) and HelpSteer2 (Wang et al., 2024b). From each dataset, we sample 500 examples, embed their contexts using OpenAI’s text-embedding-3-small model (OpenAI, 2024), and compute average pairwise cosine distances. WildChat users show the greatest contextual diversity (0.865), compared to 0.751 for HH-RLHF and 0.848 for HelpSteer2. These results suggest that real user interactions not only reflect authentic everyday needs but also span broader contexts and requests. Additional visualizations are provided in Appendix B.

**User personas are diverse with distinct characteristics.** We restructure the dataset by user and construct natural-language *personas* that summarize each individual’s preferences from their conversation histories (see prompt in Figure 7). We observe that: (1) some users provide little feedback, while others reveal clear and consistent behaviors; (2) many personas reflect common expectations, yet a notable subset exhibit unique preferences (e.g., repeatedly requesting analogies or engaging in role-play with recurring characters); and (3) some users needs vary across domains (e.g., preferring step-by-step reasoning in math but quick takeaways in daily advice) or show evolving needs over time. To study these patterns, we examine several of the most frequently mentioned preference dimensions: expertise, desired informativeness, tone, and response structure. As shown in Table 1, majorities tend to prefer expert, expansive, serious, and well-structured responses, yet substantial portions favor the opposite qualities, underscoring the need to model both dominant trends and less common preferences.

Table 1: User preferences across conversational dimensions, based on a random subset of 5,000 WildChat users. Percentages represent the proportion of users with a clear preference. “Pct. None” denotes the percentage of users with no clear preference.

Dimension	Preference 1	Pct.	Preference 2	Pct.	Pct. None
expertise	responses that can be easily understood by beginners	24.1%	responses with expert-level knowledge	<b>59.8%</b>	16.1%
informativeness	concise responses, without being verbose	36.0%	expansive and informative responses, without missing background information	<b>49.9%</b>	14.1%
tone	casual, friendly, and humorous responses	4.9%	serious, formal, and professional responses	<b>84.5%</b>	10.6%
structure	structured responses, with a clear and logical flow	<b>77.1%</b>	free-form responses, with a casual and conversational style	9.1%	13.8%

### 2.3 RLHI WITH USER-GUIDED REWRITES

In real-world scenarios, conversational models can generate unsatisfactory outputs—responses that are unhelpful, off-target, or misaligned with user intent. Organically, in such interactions, users frequently react by providing follow-up requests or explicit/implicit feedback (e.g., “Could you provide more details?”), signaling both dissatisfaction and expectations for improvement. Rather than reducing such feedback into coarse binary labels, we seek to exploit its rich semantic content. Leveraging feedback to help the model identify where it falls short and apply targeted updates provides a natural path toward more useful and better-aligned model behavior.

We rely on our user message classification in Section 2.2 to identify *re-attempts with feedback*, which make up 26.51% of all user messages in WildChat. In these cases, the model is prompted to revise its previous unsatisfactory response using the explicit or implicit user feedback (e.g., as in Figure 2, adding numbers and statistics when requested). The prompt we use is provided in Appendix Figure 8. This produces preference pairs where the user-guided rewrite is favored over the original output, directly reflecting user-indicated improvements.

To better ground learning in long-term user preferences, we prompt the LLM to summarize each user’s latent preferences from their conversation histories into a user persona. These personas are incorporated into preference pairs generated via user-guided rewrites during training, and dynamically updated at inference time to guide personalized generation, as shown in Figure 9. The persona distills long-context signals into a compact representation, while turn-level feedback offers immediate, response-specific supervision. Together, long-context persona modeling and local feedback

signals help the system capture user-specific expectations and styles that may differ from general preferences, linking a user’s enduring preferences to desirable outputs.

To ensure the quality of preference pairs, we filter the data using two criteria: (1) User-guided rewrites must improve upon the original. We discard any rewrites with a user-based reward (details in Section 2.4) lower than the original to avoid harmful follow-ups. (2) Overall quality must be high. We apply the filtering techniques from RIP (Yu et al., 2025), with details provided in Appendix C.2.

Formally, for each training instance  $i$  from user  $u$ , we consider the persona  $p_u$ , the multi-turn context  $x_{u,i}$ , a dispreferred original  $y_{u,i}^-$ , and a preferred rewrite  $y_{u,i}^+$ . We perform preference optimization using persona-conditioned Direct Preference Optimization (DPO), which maximizes the relative preference for  $y_{u,i}^+$  over  $y_{u,i}^-$  conditioned on both the prompt and persona:

$$\mathcal{L}_{\text{persona-DPO}} = \mathbb{E}_{u,i} \left[ \log \sigma \left( \beta \left( \log \frac{\pi_{\theta}(y_{u,i}^+ | x_{u,i}, p_u)}{\pi_{\text{ref}}(y_{u,i}^+ | x_{u,i}, p_u)} - \log \frac{\pi_{\theta}(y_{u,i}^- | x_{u,i}, p_u)}{\pi_{\text{ref}}(y_{u,i}^- | x_{u,i}, p_u)} \right) \right) \right], \quad (1)$$

where  $\pi_{\theta}$  is the current policy,  $\pi_{\text{ref}}$  a frozen reference model (a copy of the base model used as a baseline), and  $\beta$  controls the sharpness of preference learning. This objective explicitly conditions preference optimization on user personas, aligning generation with individualized expectations derived from long-term interactions, and yielding more personalized, satisfactory responses.

## 2.4 RLHI WITH USER-BASED REWARDS

In real-world human-LLM interactions, many initial requests do not come with follow-ups or feedback clarifying expectations for improvement. Nevertheless, these requests still reflect genuine user needs and are grounded in authentic human personas. Our goal is to improve model responses for such cases in a personalized manner. Using a (user-based) reward model provides a scalable way to learn from one-shot requests, enabling adaptation even when explicit feedback is absent.

To this end, we develop user-based rewards to guide model learning. For each user request, we generate preference pairs by first sampling  $N$  candidate responses, then evaluating them with a reward model that explicitly conditions on the corresponding user persona. For example, as illustrated in Figure 2 (right), if long-term interactions indicate that a user favors answers with numbers, statistics, and concrete evidence, the reward model will assign higher scores to responses that not only meet general quality criteria but also reflect these user-specific characteristics.

Formally, for each training instance  $i$  from user  $u$ , let  $p_u$  denote the user persona and  $x_{u,i}$  the multi-turn context. The LLM  $\mathcal{M}$  generates  $N$  candidate responses conditioned on both context and persona. A reward model  $r$  then scores each candidate given  $(x_{u,i}, p_u)$ . Preference pairs  $(y_{u,i}^+, y_{u,i}^-)$  are formed by selecting the highest- and lowest-scoring candidates:

$$\{y_{u,i}^{(n)}\}_{n=1}^N \sim \mathcal{M}(x_{u,i}, p_u) \quad \text{then} \quad \begin{cases} y_{u,i}^+ = \arg \max_{n \in [N]} r(y_{u,i}^{(n)} | x_{u,i}, p_u), \\ y_{u,i}^- = \arg \min_{n \in [N]} r(y_{u,i}^{(n)} | x_{u,i}, p_u). \end{cases} \quad (2)$$

We then apply persona-conditioned preference optimization, maximizing the relative preference for  $y_{u,i}^+$  over  $y_{u,i}^-$  given both the prompt and the persona. This can be instantiated as either offline DPO, where preference pairs are pre-collected, or online DPO, where new candidates are generated dynamically and preferences are updated on the fly. Both variants ensure that optimization is explicitly grounded in user personas, thereby complementing user-guided rewrites (Section 2.3) by extending alignment to the broader set of initial user requests when follow-up feedback is unavailable.

## 3 EXPERIMENTAL SETUP

### 3.1 TRAINING DATA GENERATION

**User Evaluation and Instruction-Following Tasks.** We build on the WildChat dataset, using 80% for training and reserving the rest for evaluation. To ensure quality, we exclude Midjourney-related

instructions and retain only users with sufficient conversation history and meaningful feedback (details in Appendix C.1). To avoid training on GPT outputs as we use Llama for training, we construct a derived dataset, *WildLlamaChat*, which preserves only user messages. Assistant responses are reconstructed by prompting Llama-3.1-8B-Instruct with the surrounding context, with details provided in Appendix C.3. For RLHI methods: (1) *RLHI with User-Guided Rewrites* uses Llama-3.1-8B-Instruct to generate user-based rewrites under sampling parameters  $T = 0.6$  and  $top-p = 0.9$ . (2) *RLHI with User-Based Rewards* samples  $N = 64$  responses per prompt from a curated pool of high-quality prompts using the same model and parameters, with the Athene-RM-8B reward model (Frick et al., 2024) providing user-based rewards.

**Reasoning Tasks.** Since no open-source dataset captures real human interactions in complex reasoning scenarios, we synthesize conversations by simulating users who ask math questions and point out model errors. These are based on the PRM800K dataset (Lightman et al., 2023), which includes MATH problems (Hendrycks et al., 2021), model-generated solutions, and step-level human correctness annotations. We randomly sample 10,000 erroneous solutions. In each conversation, the first turn presents a math problem, and the model replies with the dataset solution. In the second turn, the user makes comments such as “Step 3 seems incomplete or has an error” (details in Appendix C.4). Importantly, the simulated users only indicate where mistakes occur, without offering correct answers or detailed corrections, mimicking realistic user behavior. At training time, we apply *RLHI with User-Guided Rewrites* to revise unsatisfactory model outputs based on this feedback. Since the conversations are not tied to specific users, we do not incorporate user personas in this case.

### 3.2 TRAINING DETAILS

We initialize all models from Llama-3.1-8B-Instruct (Grattafiori et al., 2024). For RLHI methods: (1) *RLHI with User-Guided Rewrites* applies persona-conditioned DPO training, where we adopt a batch size of 64 and sweep over learning rates of  $5 \times 10^{-7}$  and  $1 \times 10^{-6}$ . (2) *RLHI with User-Based Rewards* uses persona-conditioned online DPO training with batch size 32, learning rate  $1 \times 10^{-6}$ , and KL penalty  $\beta = 0.01$ . For instruction-following tasks, we perform early stopping using the same validation set as in Yu et al. (2025).

### 3.3 MODELS AND BASELINES

We compare RLHI against the following baselines: (1) **RL with Rewrites from Scratch**, which mirrors the *RLHI with User-Guided Rewrites* pipeline, but the model regenerates its responses without access to prior outputs or user feedback; (2) **RL with User-Agnostic Rewards**, which performs online DPO training on the same prompts used in *RLHI with User-Based Rewards*, but uses generic rewards that do not consider user personas; (3) **SFT with User-Guided Rewrites** and **SFT with User-Based Rewards**, which apply supervised finetuning on the chosen responses from our generated preference pairs; and (4) **RLHI w/o Quality Filtering**, which performs *RLHI with User-Guided Rewrites* but omits quality filtering of the rewrites.

### 3.4 EVALUATION SETTING

**User-Based Evaluation.** We introduce WILDCHAT USEREVAL, an LLM-based automated evaluation of personalization and instruction-following on real-world queries. We sample 100 users from the WildChat dataset with at least 10 conversations and substantial feedback. For each user, all but the last five conversations form the reference history, and the final five multi-turn dialogues are held out for evaluation. At each user turn in the held-out set, the evaluated model generates a response, which an OpenAI o3-based judge compares against the original ChatGPT response along three axes: (1) *Personalization*, where the judge first summarizes the user’s persona from the reference history and decides which response better aligns with it; (2) *Instruction-Following*, assessing which response more faithfully follows the user’s request and provides higher-quality content; and (3) *UserEval*, a holistic judgment simulating how a user would rate the responses, incorporating both aspects (1) and (2). See Appendix G for evaluation prompts. Model outputs are generated using decoding parameters  $T = 0.6$  and  $top-p = 0.9$  (consistent across evaluations below).

We consider two inference settings: (1) *Context-Only Inference*, where the model answers using only the ongoing multi-turn context, and (2) *Persona-Guided Inference*, where the evaluated model

Table 2: **User-Based Evaluations.** Win rates (%) judged by o3 against original ChatGPT responses on WILDCHAT USEREVAL. RLHI methods achieve substantial gains in personalization, instruction-following, and overall user preference compared to the seed model and baselines.

	Personalization	Instr-Following	UserEval
<i>Baselines</i>			
Llama-3.1-8B-Instruct	38.2	30.6	32.5
+ <i>Persona-Guided Inference</i>	39.8	29.2	31.3
RL with Rewrites from Scratch	52.5	41.3	46.3
+ <i>Persona-Guided Inference</i>	54.6	40.4	47.3
RL with User-Agnostic Rewards	52.7	43.3	47.9
+ <i>Persona-Guided Inference</i>	54.2	42.8	48.4
<i>RLHI</i>			
User-Guided Rewrites	54.6	45.5	52.0
+ <i>Persona-Guided Inference</i>	<b>62.5</b>	44.5	<b>54.9</b>
User-Based Rewards	61.0	<b>46.8</b>	51.3
+ <i>Persona-Guided Inference</i>	62.3	44.7	52.5

Table 3: **Standard Evaluations.** Win rates (%) judged by GPT-4 Turbo on AlpacaEval2 and Arena-Hard. RLHI methods deliver large improvements over the seed model and baselines. *User-Based Rewards* beats or matches *RL with User-Agnostic Rewards* in this user-free setting.

<i>Standard models</i>	AlpacaEval2		Arena-Hard
	LC Win	Win	Score
Llama-3.1-8B-Instruct	20.9	21.8	21.3
RL with Rewrites from Scratch	34.7	31.0	50.0
RL with User-Agnostic Rewards	77.0	73.3	<b>64.4</b>
RLHI with User-Guided Rewrites	35.2	38.5	51.2
RLHI with User-Based Rewards	<b>77.9</b>	<b>83.4</b>	64.3

derives a persona from the reference history, and this persona is prepended to the user prompt, testing whether the model can both infer and leverage an explicit persona during generation.

To verify the reliability of LLM-based judgments, we also conduct a human study. We recruit  $N = 10$  participants, each evaluating 50 randomly sampled turns under the same *UserEval* setting, with anonymized model identities and randomized response orders. Details are provided in Appendix E.

**Standard Evaluation.** We evaluate models on AlpacaEval 2.0 (Li et al., 2023; Dubois et al., 2024) and Arena-Hard (Li et al., 2024a), which are robust instruction following benchmarks that have a high correlation with human preferences. Evaluations are conducted with GPT-4 Turbo as the judge. AlpacaEval 2.0 includes both raw and length-controlled (LC) win rates.

**Reasoning Benchmarks.** We evaluate on OlympiadBench (He et al., 2024), Minerva (Lewkowycz et al., 2022), GPQA (Rein et al., 2024), and MMLU-Pro (Wang et al., 2024a), covering diverse reasoning challenges. For each problem, we sample  $N = 50$  solutions and report average accuracy.

## 4 RESULTS

**User-Based Evaluation.** Table 2 provides results on WILDCHAT USEREVAL. RLHI methods consistently deliver strong improvements and outperform the baselines: *RLHI with User-Guided Rewrites* achieves the largest gains in personalization (+24.3) and overall improvement (+22.4), while *RLHI with User-Based Rewards* yields the strongest increase in instruction-following (+14.1). *RL with User-Agnostic Rewards* also significantly improves instruction-following but falls far behind RLHI in personalization (-8.3). Persona-guided inference enhances personalization, though sometimes at the cost of instruction-following. In the human study, *RLHI with User-Guided Rewrites* and *RLHI with User-Based Rewards* achieve win rates of 72.6% and 74.0% over Llama-3.1-8B-Instruct, confirming their effectiveness under direct human judgment.

**Standard Evaluation.** As shown in Table 3, RLHI achieves strong results in the standard user-free setting as well. *RLHI with User-Guided Rewrites* delivers large gains over Llama-3.1-8B-Instruct and outperforms *RL with Rewrites from Scratch*, although it lags behind online methods using reward models. This gap is likely due to the difference between training on multi-turn, real-user queries from WildChat and the single-turn, challenging prompts emphasized in these benchmarks. However, *RLHI with User-Based Rewards* achieves 77.9% length-controlled win rate on AlpacaEval 2.0, outperforming *RL with User-Agnostic Rewards* and ranking above all RLHF methods on the leaderboard, and matches *RL with User-Agnostic Rewards* on ArenaHard in this user-free setting.

**Reasoning Benchmarks.** As shown in Table 3, *RLHI with User-Guided Rewrites* raises average accuracy from 26.5 to 31.8 across the four reasoning benchmarks. Among them, Minerva and OlympiadBench test math reasoning, while GPQA and MMLU-Pro evaluate advanced scientific and general-domain reasoning. Although training involves only math conversations, the gains transfer beyond math to broader reasoning tasks, indicating strong generalization. Notably, unlike methods that rely on verifiable rewards or detailed annotations, our setup involves simulated users who only flag mistakes without providing correct answers or fixes. Even such lightweight, realistic feedback improves reasoning, highlighting the effectiveness of learning from natural human interaction.

Table 4: **Performance on Reasoning Benchmarks.** *RLHI with User-Guided Rewrites* consistently improves over Llama-3.1-8B-Instruct across all tasks, yielding a +5.3 average gain.

	Minerva	Olympiad	GPQA	MMLU-Pro	Avg.
Llama-3.1-8B-Instruct	20.2	14.5	26.3	44.9	26.5
RLHI with User-Guided Rewrites	<b>25.4</b>	<b>18.4</b>	<b>33.1</b>	<b>50.1</b>	<b>31.8</b>

#### 4.1 UNDERSTANDING HUMAN INTERACTION AND RLHI

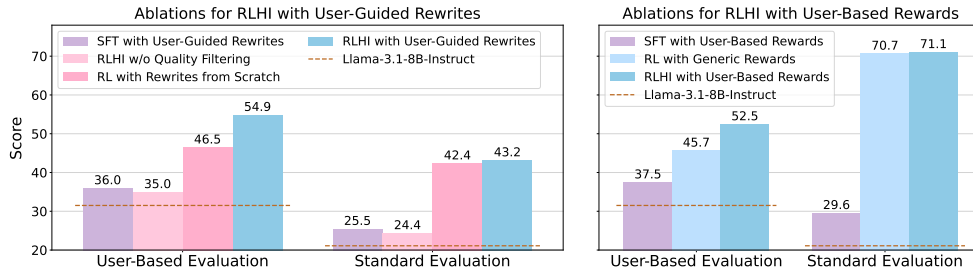


Figure 3: **Ablation Results.** User-Based Evaluation reports win rates on WILDCHAT USEREVAL, while Standard Evaluation averages AlpacaEval2 LC win rates and Arena-Hard scores. Both *RLHI with User-Guided Rewrites* and *RLHI with User-Based rewards* consistently outperform baselines.

**User-guided rewrites outperform regenerations by leveraging contextual feedback.** We compare *RLHI with User-Guided Rewrites* against *RL with Rewrites from Scratch*. When unsatisfactory responses are revised with user guidance rather than regenerated from scratch, the model benefits from direct, context-sensitive feedback that preserves the user’s original intent while correcting specific deficiencies. This leads to stronger performance, as shown by (i) head-to-head rewrite comparisons, where User-Guided Rewrites achieves a 60.4% win rate under Athene-RM-8B, and (ii) training outcomes shown in Tables 2, 3, and Figure 3, where models trained with User-Guided Rewrites outperform repeated sampling on both user-based and standard evaluations, with notably larger gains in personalization (+7.9 points).

**User-based rewards capture long-term preferences for stronger alignment.** In *RLHI with User-Based Rewards*, the reward model ranks and selects responses conditioned on a persona derived from each user’s long-term interaction history. By modeling such long-term preferences, user-based rewards guide the policy toward personalized behaviors that generalize across diverse queries. Compared to user-agnostic rewards, as shown in Tables 2, 3, and Figure 3, they substantially enhance personalization (+8.3 points), improve instruction-following and overall performance on real-world queries, and maintain competitive performance on standard benchmarks.

432 **RL outperforms supervised finetuning in learning from human interaction.** Figure 3 shows that  
 433 SFT underperforms RL across both variants of our method and both evaluations. This gap arises  
 434 because SFT relies only on positive examples and lacks gradient signals to distinguish good from  
 435 bad responses. In contrast, RL methods such as DPO optimize policies over preference signals by  
 436 leveraging both preferred and dispreferred examples, offering richer supervision regarding relative  
 437 quality and more effectively aligning models with nuanced human preferences.

438 **Human interaction data is noisy and needs quality filtering.** The main challenge in RLHI is  
 439 the noisiness of interaction data, which often includes low-quality prompts, harmful feedback,  
 440 feedback inconsistent with earlier requests, or signals misaligned with common expectations.  
 441 As shown in Figure 3, without filtering high-quality  
 442 signals using reward models, *RLHI with User-Guided*  
 443 *Rewrites* achieves only marginal gains of +2.5 and +3.3  
 444 points on user-based and standard evaluations. In contrast, filtering with reward models produces substantial  
 445 improvements of +23.4 and +17.7 points, underscoring  
 446 the critical role of quality control in leveraging human  
 447 interaction for alignment.

448 **RLHI benefits from user diversity.** *RLHI with User-*  
 449 *Guided Rewrites* learns from user conversations span-  
 450 ning 1268 users, each contributing only a few inter-  
 451 actions. To isolate the role of diversity, we construct  
 452 equally sized datasets but drawn from just 10 users  
 453 with many conversations each. As shown in Figure 4,  
 454 broader user diversity consistently improves win rates  
 455 and scales more effectively, as the model learns to adapt  
 456 to a wider range of preferences and interaction styles.

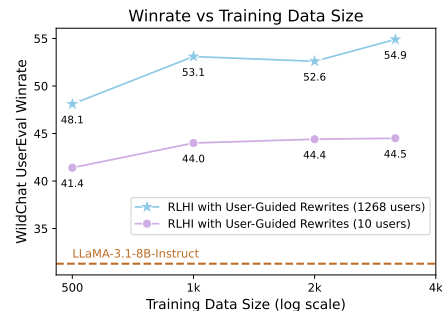


Figure 4: **Effect of user diversity on RLHI.** Training with 1268 diverse users outperforms training with 10 users of similar data size on WILDCHAT USEREVAL.

## 459 5 RELATED WORK

462 **Learning from Human Feedback.** Reinforcement Learning from Human Feedback (RLHF) trains  
 463 a reward model on preference data and optimizes the base model with RL (Ziegler et al., 2019; Sti-  
 464 ennon et al., 2020; Ouyang et al., 2022). Later work replaces explicit RL with direct preference opti-  
 465 mization and related objectives for greater stability and efficiency (Rafailov et al., 2023; Ethayarajh  
 466 et al., 2024; Azar et al., 2024). Beyond curated datasets, feedback is increasingly mined from post-  
 467 deployment interactions: using user message classifiers (Hancock et al., 2019; Chen et al., 2024b;  
 468 Don-Yehiya et al., 2024; Han et al., 2025), heuristics such as response length (Pang et al., 2023), or  
 469 organic user signals like thumbs up/down and free-form comments (Jaques et al., 2020; Xu et al.,  
 470 2023) to assess user attitude or satisfaction. These signals are then optimized via fine-tuning (Don-  
 471 Yehiya et al., 2024) or other methods (Xu et al., 2023; Pang et al., 2023). Unlike prior work that  
 472 relies on annotated labels or proxy signals, WildFeedback (Shi et al., 2024) and our RLHI approach  
 473 learn directly from organic interactions. Our work goes further by modeling long-term user history,  
 474 proposing user-based rewards, demonstrating stronger performance on standard benchmarks, and  
 475 enabling user-based evaluation.

476 **Personalizing Language Models.** Personalization aims to adapt LMs to user preferences through  
 477 retrieval, prompting, representation learning, or RLHF (Zhang et al., 2024). Retrieval and prompt-  
 478 ing approaches incorporate user information as external memory (Mysore et al., 2023; Salemi et al.,  
 479 2024) or as persona/profile context (Jiang et al., 2023). Representation-learning methods encode  
 480 traits in model parameters (Tan et al., 2024) or embeddings (Chen et al., 2025). RLHF-style per-  
 481 sonalization uses user information as reward signals to align LLMs with personalized preferences:  
 482 works explore conditioning on multiple reward dimensions (Jang et al., 2023; Yang et al., 2024; Li  
 483 et al., 2024b; Shenfeld et al., 2025), decoupling generation dynamics from user utility (Chen et al.,  
 484 2024a), generalized system messages during training (Lee et al., 2024), or aligning models through  
 485 a user-specific latent variable model (Poddar et al., 2024). Our RLHI framework explicitly con-  
 nects long-term personas with turn-level preferences and optimizes on organic interactions, yielding  
 stronger personalization and better instruction-following.

---

## 6 CONCLUSION

In this paper, we make the case for the improvement of models by learning from real-world human interaction. We present a concrete method, Reinforcement Learning from Human Interaction (RLHI), a simple and scalable framework for learning directly from in-the-wild user conversations utilizing long-term conversation history and organic natural-language feedback. RLHI provides clear improvements when measured at the *user* level compared to strong baselines, where utilizing organic feedback is shown to improve both non-reasoning and reasoning tasks. Looking forward, we see opportunities to extend RLHI with human-in-the-loop learning, richer and safer reward modeling, privacy-preserving personalization, and broader modality and task coverage. Importantly, we believe using RLHI within an online learning loop, where a continually updating deployed model learns from its organic interactions, would bring major gains compared to the fixed training data setup in our experiments. We hope these findings encourage a shift toward learning from real-world human interaction to build capable, personalized assistants that improve over time.

## REFERENCES

- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR, 2024.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai, and Zuozhu Liu. Pad: Personalized alignment of llms at decoding-time. *arXiv preprint arXiv:2410.04070*, 2024a.
- Runjin Chen, Andy Arditi, Henry Sleight, Owain Evans, and Jack Lindsey. Persona vectors: Monitoring and controlling character traits in language models. *arXiv preprint arXiv:2507.21509*, 2025.
- Zizhao Chen, Mustafa Omer Gul, Yiwei Chen, Gloria Geng, Anne Wu, and Yoav Artzi. Retrospective learning from interactions. *arXiv preprint arXiv:2410.13852*, 2024b.
- Shachar Don-Yehiya, Leshem Choshen, and Omri Abend. Naturally occurring feedback is common, extractable and useful. *arXiv preprint arXiv:2407.10944*, 2024.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- Evan Frick, Peter Jin, Tianle Li, Karthik Ganesan, Jian Zhang, Jiantao Jiao, and Banghua Zhu. Athene-70b: Redefining the boundaries of post-training for open models, july 2024. URL <https://huggingface.co/Nexusflow/Athene-70B>, 2024.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Eric Han, Jun Chen, Karthik Abinav Sankararaman, Xiaoliang Peng, Tengyu Xu, Eryk Helenowski, Kaiyan Peng, Mrinal Kumar, Sinong Wang, Han Fang, et al. Reinforcement learning from user feedback. *arXiv preprint arXiv:2505.14946*, 2025.
- Braden Hancock, Antoine Bordes, Pierre-Emmanuel Mazare, and Jason Weston. Learning from dialogue after deployment: Feed yourself, chatbot! *arXiv preprint arXiv:1901.05415*, 2019.

---

540 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,  
541 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for  
542 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint*  
543 *arXiv:2402.14008*, 2024.

544 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
545 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*  
546 *preprint arXiv:2103.03874*, 2021.

548 Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer,  
549 Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Per-  
550 sonalized large language model alignment via post-hoc parameter merging. *arXiv preprint*  
551 *arXiv:2310.11564*, 2023.

552 Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharioun, Craig Ferguson, Agata Lapedriza, Noah  
553 Jones, Shixiang Shane Gu, and Rosalind Picard. Human-centric dialog training via offline rein-  
554 forcement learning. *arXiv preprint arXiv:2010.05848*, 2020.

556 Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. Evalu-  
557 ating and inducing personality in pre-trained language models. *Advances in Neural Information*  
558 *Processing Systems*, 36:10622–10643, 2023.

560 Seongyun Lee, Sue Hyun Park, Seungone Kim, and Minjoon Seo. Aligning to thousands of prefer-  
561 ences via system message generalization. *Advances in Neural Information Processing Systems*,  
562 37:73783–73829, 2024.

563 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-  
564 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative  
565 reasoning problems with language models. *Advances in neural information processing systems*,  
566 35:3843–3857, 2022.

568 Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gon-  
569 zalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and  
570 benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024a.

571 Xinyu Li, Ruiyang Zhou, Zachary C Lipton, and Liu Leqi. Personalized language modeling from  
572 personalized human feedback. *arXiv preprint arXiv:2402.05133*, 2024b.

574 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy  
575 Liang, and Tatsunori B Hashimoto. Alpacaeval: An automatic evaluator of instruction-following  
576 models, 2023.

577 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan  
578 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*  
579 *International Conference on Learning Representations*, 2023.

581 Sheshera Mysore, Zhuoran Lu, Mengting Wan, Longqi Yang, Steve Menezes, Tina Baghaee, Em-  
582 manuel Barajas Gonzalez, Jennifer Neville, and Tara Safavi. Pearl: Personalizing large language  
583 model writing assistants with generation-calibrated retrievers. *arXiv preprint arXiv:2311.09180*,  
584 2023.

585 OpenAI. text-embedding-3-small, 2024. URL [https://platform.openai.com/docs/  
586 guides/embeddings](https://platform.openai.com/docs/guides/embeddings).

588 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
589 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
590 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
591 27730–27744, 2022.

592 Richard Yuanzhe Pang, Stephen Roller, Kyunghyun Cho, He He, and Jason Weston. Leveraging  
593 implicit feedback from deployment data in dialogue. *arXiv preprint arXiv:2307.14117*, 2023.

---

594 Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. Personalizing  
595 reinforcement learning from human feedback with variational preference learning. *Advances in*  
596 *Neural Information Processing Systems*, 37:52516–52544, 2024.

597

598 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
599 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*  
600 *in neural information processing systems*, 36:53728–53741, 2023.

601

602 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-  
603 rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-  
604 mark. In *First Conference on Language Modeling*, 2024.

605

606 Alireza Salemi, Surya Kallumadi, and Hamed Zamani. Optimization methods for personalizing  
607 large language models through retrieval augmentation. In *Proceedings of the 47th International*  
608 *ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 752–762,  
609 2024.

610

611 Idan Shenfeld, Felix Faltings, Pulkit Agrawal, and Aldo Pacchiano. Language model personalization  
612 via reward factorization. *arXiv preprint arXiv:2503.06358*, 2025.

613

614 Taiwei Shi, Zhuoer Wang, Longqi Yang, Ying-Chun Lin, Zexue He, Mengting Wan, Pei Zhou, Sujay  
615 Jauhar, Sihao Chen, Shan Xia, et al. Wildfeedback: Aligning llms with in-situ user interactions  
616 and feedback. *arXiv preprint arXiv:2408.15549*, 2024.

617

618 David Silver and Richard S Sutton. Welcome to the era of experience. *Google AI*, 1, 2025.

619

620 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
621 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances*  
622 *in neural information processing systems*, 33:3008–3021, 2020.

623

624 Zhaoxuan Tan, Zheyuan Liu, and Meng Jiang. Personalized pieces: Efficient personalized large  
625 language models through collaborative efforts. *arXiv preprint arXiv:2406.10471*, 2024.

626

627 Michael Tomasello, Malinda Carpenter, Josep Call, Tanya Behne, and Henrike Moll. Understanding  
628 and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences*, 28(5):  
629 675–691, 2005.

630

631 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming  
632 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-  
633 task language understanding benchmark. *Advances in Neural Information Processing Systems*,  
634 37:95266–95290, 2024a.

635

636 Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy Zhang,  
637 Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. Helpsteer 2: Open-source dataset for training  
638 top-performing reward models. *Advances in Neural Information Processing Systems*, 37:1474–  
639 1501, 2024b.

640

641 Jing Xu, Da Ju, Joshua Lane, Mojtaba Komeili, Eric Michael Smith, Megan Ung, Morteza Behrooz,  
642 William Ngan, Rashed Moritz, Sainbayar Sukhbaatar, et al. Improving open language models by  
643 learning from organic interactions. *arXiv preprint arXiv:2306.04707*, 2023.

644

645 Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. Rewards-  
646 in-context: Multi-objective alignment of foundation models with dynamic preference adjustment.  
647 *arXiv preprint arXiv:2402.10207*, 2024.

648

649 Ping Yu, Weizhe Yuan, Olga Golovneva, Tianhao Wu, Sainbayar Sukhbaatar, Jason Weston, and  
650 Jing Xu. Rip: Better models by survival of the fittest prompts. *arXiv preprint arXiv:2501.18578*,  
651 2025.

652

653 Weizhe Yuan, Ilia Kulikov, Ping Yu, Kyunghyun Cho, Sainbayar Sukhbaatar, Jason Weston, and  
654 Jing Xu. Following length constraints in instructions. *arXiv preprint arXiv:2406.17744*, 2024.

648 Zhehao Zhang, Ryan A Rossi, Branislav Kveton, Yijia Shao, Diyi Yang, Hamed Zamani, Franck  
649 Dernoncourt, Joe Barrow, Tong Yu, Sungchul Kim, et al. Personalization of large language mod-  
650 els: A survey. *arXiv preprint arXiv:2411.00027*, 2024.  
651  
652 Hao Zhao, Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Long is more  
653 for alignment: A simple but tough-to-beat baseline for instruction fine-tuning. *arXiv preprint*  
654 *arXiv:2402.04833*, 2024a.  
655  
656 Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat:  
657 1m chatgpt interaction logs in the wild. *arXiv preprint arXiv:2405.01470*, 2024b.  
658  
659 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
660 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*  
661 *preprint arXiv:1909.08593*, 2019.  
662

## 664 A EXAMPLES OF USER CONVERSATION DATA

665  
666 As described in Section 2.2, each multi-turn conversation begins with an *initial request*, which  
667 accounts for 27.07% of all user messages. User follow-ups are then classified into four types: *new*  
668 *requests* (40.40%), *re-attempts with feedback* (26.51%), *re-attempts without feedback* (4.77%), and  
669 *positive feedback* (1.25%). Below, we provide real examples from WildChat for each type.  
670

### 671 A.1 NEW REQUESTS

672  
673 We present three examples of *new requests*: (1) In the first, the *initial request* is only a greeting  
674 or non-question, followed by a substantive query. (2) In the second, the user makes a different  
675 translation request. Although it is still thematically related to the initial one (both about Viennese  
676 translation), the specific content changes, so it is classified as a *new request*. (3) In the third, the user  
677 shifts topics entirely, making a request unrelated to the original.  
678

#### 679 NEW REQUEST – EXAMPLE 1

680 **User message:** hi

681 **Assistant message:** Hello! How can I assist you today?

682 **User message:** write a long poem about AI

#### 686 NEW REQUEST – EXAMPLE 2

687 **User message:** Hello! Please tell me how to say “I am also a kid” in Viennese Dialect

688 **Assistant message:** In Viennese dialect, you would say “I bin aa a Kind.”

689 **User message:** So, how is your day? In Viennese

#### 694 NEW REQUEST – EXAMPLE 3

695 **User message:** how begin conversation with girl that i like

696 **Assistant message:** As an AI language model, I do not have personal experiences or emo-  
697 tions, but here are some tips on how to start a conversation with someone you like:  
698 1. Start with... (713 characters skipped here)

699 **User message:** teach me speak english  
700  
701

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A.2 RE-ATTEMPTS WITH FEEDBACK

We present three examples of *re-attempts with feedback*: (1) In the first, the user repeats the request for more ideas but adds the instruction “make them creative,” offering explicit guidance for improvement. (2) In the second, the user reframes the request by introducing a new dimension, asking how social classes perceive and interpret an act, which meaningfully extends the original query. This reveals the perspective the user values and signals what could be incorporated directly into future responses. (3) In the third, the user corrects the assistant by pointing out a length requirement (“2000+ words”), providing clear feedback on where the response fell short and how to adjust. Interestingly, the initial request specified “a 2000 word essay,” without mentioning “2000+” or “about 2000.” This inconsistency may reflect a linguistic habit, but it also helps the assistant better interpret the user’s intended word limit in future interactions.

### RE-ATTEMPT WITH FEEDBACK – EXAMPLE 1

**User message:** give me 20 ideas for themes for a summer camp for children aging from 4 to 18

**Assistant message:** 1. Nature and Outdoor Adventure... (642 characters skipped here)

**User message:** give me 20 more make them creative

### RE-ATTEMPT WITH FEEDBACK – EXAMPLE 2

**User message:** What cultural or traditional significance does wearing a hood hold in this particular world?

**Assistant message:** In this particular world, wearing a hood holds cultural and traditional significance in several ways:

1. Protection from the elements... (1946 characters skipped here)

**User message:** How do different social classes in this world perceive and interpret the act of wearing a hood?

### RE-ATTEMPT WITH FEEDBACK – EXAMPLE 3

**User message:** write a 2000 word essay on how political ideologies influence modern political thinking

**Assistant message:** Introduction:

Political ideologies have been around for centuries, and they have influenced the way people think about politics and governance... (9406 characters skipped here)

**User message:** i said 2000+ word

## A.3 RE-ATTEMPTS WITHOUT FEEDBACK

We present two examples of *re-attempts without feedback*. (1) In the first, the user repeats the same request but merely adds “Best,” signaling dissatisfaction without offering actionable guidance. (2) In the second, the user initially asks for 20 more examples with the instruction “make them creative,” but later repeats the request for 20 more without providing any new signals, which is then classified as a *re-attempt without feedback*.

### RE-ATTEMPT WITHOUT FEEDBACK – EXAMPLE 1

**User message:** give me some prompt for ask Ai to design to me wallpeper phone

**Assistant message:** Sure, here are some creative prompts you can use to ask an AI to design a wallpaper for your phone:... (1698 characters skipped here)

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

**User message:** give me some Best prompt for ask Ai to design to me wallpeper phone

**RE-ATTEMPT WITHOUT FEEDBACK – EXAMPLE 2**

**User message:** give me 20 ideas for themes for a summer camp for children aging from 4 to 18  
**Assistant message:** 1. Nature and Outdoor Adventure... (642 characters skipped here)  
**User message:** give me 20 more make them creative  
**Assistant message:** 1. Fantasy Island Adventure... (702 characters skipped here)  
**User message:** give me 20 more

**A.4 POSITIVE FEEDBACK**

Positive feedback is rare, accounting for only 1.25% of user messages, since users often end the conversation once satisfied. Below, we show one example where the user expresses appreciation with a simple “thank you.”

**POSITIVE FEEDBACK – EXAMPLE 1**

**User message:** What is the difference between Hinduism and Buddhism ?  
**Assistant message:** As an AI language model, I don’t believe in beliefs and religions. However, here are the major differences between Hinduism and Buddhism:  
1. Belief in God:... (1267 characters skipped here)  
**User message:** Thank you  
**Assistant message:** You’re welcome!

**B DETAILS OF USER CONVERSATION DATA ANALYSIS**

We visualize the prompt embeddings calculated in Section 2.2. This demonstrates that real-world human interaction data are more diverse than existing preference datasets.

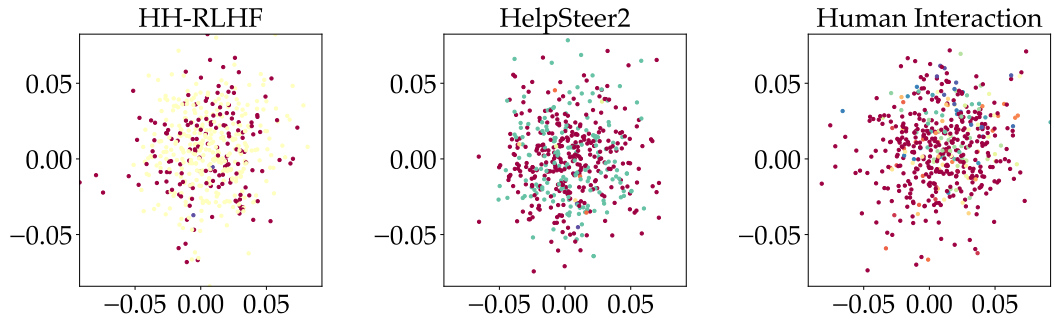


Figure 5: Visualization of context embeddings across preference datasets: the two annotated human feedback datasets, HH-RLHF and HelpSteer2, and our human interaction dataset used for RLHI.

---

## 810 C DETAILS OF USER CONVERSATION DATA PROCESSING

### 811 C.1 DETAILS OF TRAINING CONVERSATION FILTERING

812 To ensure data quality and relevance, we apply several filtering steps to the WildChat-1M dataset  
813 (Zhao et al., 2024b), before using RLHI to learn from the user conversations:

- 814 1. Exclude non-English prompts using the provided language annotations.
- 815 2. Remove Midjourney-related instructions, which typically begin with: “As a prompt gener-  
816 ator for a generative AI called ‘Midjourney’, you will create image prompts ...”.
- 817 3. Retain only users with at least three conversations, ensuring enough context to infer a per-  
818 sona.
- 819 4. Discard users with more than 100 conversations, as they are often associated with program-  
820 generated instructions that are low quality and misaligned with real human needs.
- 821 5. Exclude conversations with more than 10 turns to maintain task focus and coherence.
- 822 6. Use an LLM to filter for users who provide meaningful feedback.

### 823 C.2 DETAILS OF PREFERENCE PAIR FILTERING

824 To improve the quality of preference pairs used for optimization, we adopt RIP’s filtering techniques  
825 (Yu et al., 2025) with the following thresholds:

- 826 1. **Rejected response length**  $\geq 1878$ : Following Yu et al. (2025), we treat rejected response  
827 length as a proxy for prompt quality. Low-quality prompts (unclear, ambiguous, or con-  
828 flicting) tend to produce short, uninformative responses, which correlate with weaker per-  
829 formance (Zhao et al., 2024a; Yuan et al., 2024).
- 830 2. **Rejected response reward**  $\geq -1$ : We use Athene-RM-8B (Frick et al., 2024) to assign  
831 user-based rewards, ensuring rejected responses still meet a minimal quality threshold.
- 832 3. **Reward gap**  $\leq 1$ : Large reward gaps often arise from low-quality prompts that allow  
833 multiple interpretations. By restricting the gap between chosen and rejected responses, we  
834 favor prompts that elicit consistent, high-quality outputs.

### 835 C.3 DETAILS OF CONSTRUCTING WILDLLAMACHAT

836 To avoid training on GPT outputs as we use Llama for training, we construct a derived dataset,  
837 *WildLlamaChat*, which preserves only user messages. Assistant responses are reconstructed by  
838 prompting Llama-3.1-8B-Instruct with the surrounding context. We don’t just use the previous con-  
839 text—we also provide the subsequent user messages. This is crucial because the follow-up feedback  
840 naturally constrains the reconstruction to be consistent with that feedback. For example, if a user’s  
841 next message says “rewrite it, do not consider it from the angle of A,” the reconstructed response  
842 will be generated to consider angle A, making it appropriate for the user’s feedback. We empiri-  
843 cally found that reconstruction with only previous context does produce misaligned responses, but  
844 including future context significantly improves alignment. Additionally, we apply reward model-  
845 based filtering to ensure quality and relevance. While we acknowledge this approach has limitations,  
846 it’s necessary to avoid training on GPT outputs while still leveraging valuable user feedback from  
847 WildChat.

848 We use Llama-3.1-8B-Instruct to both generate the training data (reconstructions and rewrites) and  
849 serve as the model being trained. This is essentially a self-improvement setup where the model  
850 learns from its own responses and user feedback on those responses. If we were to train a different  
851 model, we would use that specific model to generate its own training data, maintaining consistency  
852 between the data generation and training processes.

### 853 C.4 DETAILS OF SYNTHESIZING MATH CONVERSATIONS

854 Since no open-source dataset captures real human interactions in complex reasoning scenarios, we  
855 synthesize conversations by simulating users who ask math questions and point out model errors.

These are based on the PRM800K dataset (Lightman et al., 2023), which includes MATH problems (Hendrycks et al., 2021), model-generated solutions, and step-level human correctness annotations. From this corpus, we randomly sample 10,000 erroneous solutions and the corresponding questions.

Each synthetic conversation begins with a math problem ending with the instruction: “Please reason step by step, and put your final answer within `\boxed{\}`.” The model then replies with the dataset solution, consisting of multiple steps annotated with human judgments of correctness. In the next turn, the user identifies the first incorrect step and provides natural-language comments such as “Step 3 seems incomplete or has an error.” If the final answer is correct despite earlier mistakes, the user adds a qualifier such as “... though your final answer is correct.” In this way, the simulated users only indicate where mistakes occur, without offering correct answers or detailed corrections, mimicking realistic user behavior.

## D ADDITIONAL RESULTS ON WILDCHAT USEREVAL

In Table 5, we show results on WILDCHAT USEREVAL, breaking down the overall win rates from Table 2 into performance on initial turns and following turns. This decomposition reveals how well models handle first attempts compared to user follow-ups later in the conversation. RLHI methods continue to outperform baselines across both settings, with the strongest gains from User-Guided Rewrites, which achieves 60.3% on initial turns and 52.6% on follow-up turns when combined with Persona-Guided Inference, leading to the best overall UserEval score of 54.9%. These results highlight that RLHI consistently enhances model responses throughout multi-turn interactions.

Table 5: **User-Based Evaluations with Turn-Level Breakdown.** Win rates (%) judged by o3 against original ChatGPT responses on WILDCHAT USEREVAL. This table expands upon the UserEval results in Table 2 by separately reporting performance on initial user turns (“Initial”) and follow-up turns (“Follow-up”), providing a more detailed view of how models handle different types of requests.

	UserEval (Initial)	UserEval (Follow-up)	UserEval
Llama-3.1-8B-Instruct	36.3	30.9	32.5
+ <i>Persona-Guided Inference</i>	33.0	30.6	31.3
RL with Rewrites from Scratch	47.2	45.9	46.3
+ <i>Persona-Guided Inference</i>	46.7	47.6	47.3
RL with User-Agnostic Rewards	50.6	46.8	47.9
+ <i>Persona-Guided Inference</i>	50.6	47.5	48.4
RLHI with User-Guided Rewrites	57.0	49.9	52.0
+ <i>Persona-Guided Inference</i>	<b>60.3</b>	<b>52.6</b>	<b>54.9</b>
RLHI with User-Based Rewards	50.3	51.7	51.3
+ <i>Persona-Guided Inference</i>	54.7	51.6	52.5

## E DETAILS OF THE HUMAN STUDY

We recruited 10 participants to evaluate the model outputs, none of whom were paper authors. Participants gave informed consent under an IRB-exempt protocol and were compensated at standard rates. Each participant evaluated 50 items sampled from held-out WildChat conversations. For each item, we showed the multi-turn context up to the user’s last message and a short persona summary derived from that user’s past conversations (the same prompt used in our automated “UserEval” setup). Raters then saw two anonymized responses (A/B) from different models in random order and selected which they would prefer as the user, considering both how well it followed instructions and how personalized it was to the user’s history. We compared Llama-3.1-8B-Instruct to each RLHI variant. In total, the study produced 500 pairwise judgments: 250 for Llama-3.1-8B-Instruct vs. RLHI with User-Guided Rewrites and 250 for Llama-3.1-8B-Instruct vs. RLHI with User-Based Rewards.

918 The results of our human study strongly supported the findings from our automated evaluation.  
919 Human evaluators preferred the responses from our RLHI with User-Guided Rewrites method 72.6%  
920 of the time over the baseline model, and they preferred the RLHI with User-Based Rewards method  
921 74.0% of the time. This close alignment between human judgments and our automated metrics  
922 indicates that our methods are effective in generating responses that real users find more helpful and  
923 personalized.

## 924 F PROMPTS USED IN RLHI

925 We provide the prompts used in RLHI methods, including those for classifying user messages,  
926 inferring user personas, generating user-guided rewrites, and performing persona-guided inference.  
927

### 928 CLASSIFYING USER MESSAGES

929 You are given two requests from a user during their conversation with an AI assistant. Clas-  
930 sify the second request in relation to the first using the following labels:

931 [New] A new topic or task, or a significantly different variation of the previous task.

932 [Re-attempt with feedback] A re-attempt of the same task that includes explicit or implicit  
933 feedback, or a revised prompt.

934 [Re-attempt without feedback] A repeat of the same task, without any feedback.

935 [Positive feedback] A signal of praise or satisfaction with the previous response.

936 1st request: Write a short poem about the ocean.

937 2nd request: What's the capital of Japan?

938 Classification: [[New]]

939 1st request: Write a short poem about the ocean.

940 2nd request: Write a short poem about the ocean.

941 Classification: [[Re-attempt without feedback]]

942 1st request: Write a short poem about the ocean.

943 2nd request: Can you make it more rhyme?

944 Classification: [[Re-attempt with feedback]]

945 1st request: {initial\_request}

946 2nd request: {current\_request}

947 Classification:

948 Figure 6: Prompt for classifying user follow-up messages into four types: (1) new requests, (2)  
949 re-attempts with feedback, (3) re-attempts without feedback, and (4) positive feedback.  
950

### 951 INFERRING USER PERSONA

952 Below are user messages from conversations between this user and an AI assistant. Please  
953 list up to five key points that capture how the user prefer the assistant to respond. Output  
954 only the inferred preference, without any additional commentary or explanation.

955 [The Start of User Messages]

956 {user\_message\_history}

957 [The End of User Messages]

958 Figure 7: Prompt for deriving a natural-language user persona given each user's long-term conver-  
959 sational history.  
960

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

**GENERATING USER-GUIDED REWRITES**

Please revise your previous response based on the user feedback or follow-up request below. Ensure the revised response is not significantly longer, unless the user explicitly requests so. Ensure the revised response adheres to safety and ethical guidelines, even if the user suggests otherwise. Do not reference or mention the user feedback in your response. Output only the revised response, without any additional commentary or explanation.

[The Start of User Follow-up Response]  
{user\_response}  
[The End of User Follow-up Response]

Figure 8: Prompt for revising unsatisfactory model outputs based on users’ natural-language follow-up responses.

**SYSTEM PROMPT FOR PERSONA-GUIDED INFERENCE**

You are a helpful and personalized assistant. Prioritize your responses based on the user’s current request and conversational context. When appropriate, tailor your responses to align with the user persona provided below.  
User persona: {user\_persona}

Figure 9: System prompt for persona-guided inference. At inference time, incorporating this lightweight prompt enables the model to generate personalized responses. During training, RLHI integrates the same prompt into preference pairs, allowing the model to learn the connection between a user’s long-term persona and their turn-level, context-specific preferences.

## G PROMPTS USED IN WILDCHAT USEREVAL

We provide the prompts used in WILDCHAT USEREVAL, including those for judging personalization, instruction-following, and UserEval.

**PERSONALIZATION JUDGE**

You are given a conversation history that ends with a user question, followed by two responses from two AI assistants. You are also provided with a user persona that describes how the user prefers the assistant to respond. Your task is to act as an impartial judge and determine which response better aligns with the user persona. Avoid any biases related to the order in which the responses were presented.

Provide your verdict strictly following this format:  
- Only output “[A]” if Assistant A is better  
- Only output “[B]” if Assistant B is better

[The Start of Conversation History]  
{conversation\_history}  
[The End of Conversation History]

[The Start of Assistant A’s Answer]  
{response\_A}  
[The End of Assistant A’s Answer]

[The Start of Assistant B’s Answer]  
{response\_B}  
[The End of Assistant B’s Answer]

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

[The Start of User Persona]  
{persona}  
[The End of User Persona]

Figure 10: Prompt for the personalization judge in WILDCHAT USEREVAL. The judge first summarizes the user’s persona from the reference history using the prompt in Figure 7, and then applies this prompt to determine which response aligns better with it.

**INSTRUCTION-FOLLOWING JUDGE**

You are given a conversation history that ends with a user question, followed by two responses from two AI assistants. Your task is to act as an impartial judge and determine which response better follows the user’s instructions and provides a higher-quality answer. Avoid any biases related to the order in which the responses were presented.

Provide your verdict strictly following this format:

- Only output “[A]” if Assistant A is better
- Only output “[B]” if Assistant B is better

[The Start of Conversation History]  
{conversation\_history}  
[The End of Conversation History]

[The Start of Assistant A’s Answer]  
{response\_A}  
[The End of Assistant A’s Answer]

[The Start of Assistant B’s Answer]  
{response\_B}  
[The End of Assistant B’s Answer]

Figure 11: Prompt for the instruction-following judge in WILDCHAT USEREVAL, determining which response better follows the user’s instructions and provides a higher-quality answer.

**USEREVAL JUDGE**

You are given a conversation history that ends with a user question, followed by two responses from two AI assistants. You are also provided with a user persona that describes how the user prefers the assistant to respond. Your task is to act as an impartial judge, simulating how the user would evaluate the responses. Specifically, determine which response better follows the user’s instructions, provides a higher-quality answer, and aligns with the user persona. Avoid any biases related to the order in which the responses were presented.

Provide your verdict strictly following this format:

- Only output “[A]” if Assistant A is better
- Only output “[B]” if Assistant B is better

[The Start of Conversation History]  
{conversation\_history}  
[The End of Conversation History]

[The Start of Assistant A’s Answer]  
{response\_A}

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

```
[The End of Assistant A's Answer]

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

[The Start of User Persona]
{persona}
[The End of User Persona]
```

Figure 12: Prompt for the UserEval judge in WILDCHAT USEREVAL. The judge first summarizes the user’s persona from the reference history using the prompt in Figure 7, and then applies this prompt to determine which response better follows the user’s instructions, provides a higher-quality answer, and aligns with the user’s persona.

## H THE USE OF LARGE LANGUAGE MODELS

In accordance with the ICLR 2026 Author Guide, we disclose our use of Large Language Models (LLMs): after completing the draft of the paper, LLMs were used to polish the writing. They were not used for any other purpose.