

# ADVANCING LLM REASONING GENERALISTS WITH PREFERENCE TREES

Anonymous authors

Paper under double-blind review

## ABSTRACT

We introduce **EURUS**, a suite of large language models (LLMs) optimized for reasoning. Finetuned from Mistral-7B, Llama-3-8B, and Mixtral-8x22B, EURUS models achieve state-of-the-art results among open-source models on a diverse set of benchmarks covering mathematics, code generation, and logical reasoning problems. Notably, EURUX-8x22B outperforms GPT-3.5 Turbo in reasoning through a comprehensive benchmarking across 12 test sets covering five tasks. The strong performance of EURUS can be primarily attributed to **ULTRAINTERACT**, our newly-curated large-scale, high-quality training data dataset specifically designed for complex reasoning tasks. **ULTRAINTERACT** can be used in both supervised fine-tuning, preference learning, and reward modeling. It pairs each instruction with a preference tree consisting of (1) reasoning chains with diverse planning strategies in a unified format, (2) multi-turn interaction trajectories with the environment and the critique, and (3) pairwise positive and negative responses to facilitate preference learning. **ULTRAINTERACT** allows us to conduct an in-depth exploration of preference learning for reasoning tasks. Our investigation reveals that some well-established preference learning algorithms may be less suitable for reasoning tasks compared to their effectiveness in general conversations. The hypothesis is that in reasoning tasks, the space of correct answers is much smaller than that of incorrect ones, so it is necessary to explicitly increase the reward of chosen data. Therefore, in addition to increasing the reward margin as many preference learning algorithms do, the absolute values of positive responses' rewards should be positive and may serve as a proxy for performance. Inspired by this, we derive a novel reward modeling objective and empirically that it leads to a stable reward modeling curve and better performance. Together with **ULTRAINTERACT**, we obtain a strong reward model.

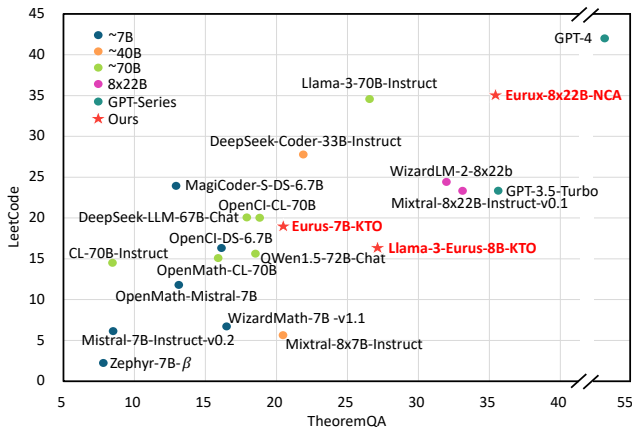


Figure 1: Evaluation results on LeetCode and TheoremQA, two challenging OOD coding and math benchmarks with only test sets. Our EURUS-7B is comparable with baselines that are 10x larger and EURUX-8x22B is the only one on par with GPT-3.5 Turbo.

## 1 INTRODUCTION

Current alignment techniques have significantly advanced the development of open-source large language models (LLMs) that effectively meet user expectations and align with human values (Touvron

054 et al., 2023; Tunstall et al., 2023). On complex reasoning, success has been achieved by specializing  
055 models for specific capabilities, such as coding (Wei et al., 2023; Guo et al., 2024a; Zheng et al., 2024)  
056 and solving math problems (Fu et al., 2023; Yue et al., 2023; Luo et al., 2023a; Toshniwal et al., 2024).  
057 However, these models still fall short, by large margins, of the most advanced proprietary models in  
058 their all-around capabilities to tackle a diverse range of challenging problems. We conjecture that this  
059 performance gap can be primarily attributed to (1) the lack of high-quality alignment data and (2)  
060 the underexploration of preference learning techniques for improving models’ complex reasoning  
061 capabilities. In this paper, we take strides towards bridging this gap by addressing both factors and  
062 developing EURUS.

063 EURUS consists of a suite of LLMs and reward model finetuned from Mistral-7B (Jiang et al.,  
064 2023a), Llama-3 (Meta, 2024), and Mixtral-8x22B (Jiang et al., 2024). Across a diverse set of  
065 complex reasoning benchmarks that are mostly out-of-distribution (OOD), EURUS achieves state-of-  
066 the-art overall performance among all open-source models. In particular, EURUS excels in solving  
067 challenging problems that often require sophisticated planning, reasoning, tool integration, and the  
068 ability to interact with and learn from the environment and users. As shown in Figure 1, on university-  
069 level STEM questions TheoremQA (Chen et al., 2023) and competition-level coding problems  
070 LeetCode Contest (Guo et al., 2024a), EURUS significantly outperforms all open-source models,  
071 achieving comparable performance to GPT-3.5 Turbo. Besides, our reward model, EURUS-RM-7B,  
072 outperforms baselines with a 5x larger size on various reward model benchmarks and demonstrates  
073 superiority on best-of-N and MCTS-guided decoding tasks.

074 EURUS models are trained on ULTRAINTERACT, our newly-curated, large-scale, and high-quality  
075 alignment data specifically designed to improve LLMs’ reasoning capabilities. ULTRAINTERACT  
076 consists of a diverse set of instructions spanning math, coding, and logical reasoning problems from 12  
077 established datasets. For each instruction, ULTRAINTERACT collects a preference tree that includes:  
078 (1) **Diverse planning strategies in a unified pattern**, such as sequential processing (Wei et al., 2022)  
079 and tool creation (Qian et al., 2023), followed by executing step-by-step actions formatted in either  
080 text or code, to provide diverse reasoning trajectories. (2) **Multi-turn interaction trajectories with  
081 the environment and the critique**, to improve models’ capabilities to learn from feedback and correct  
082 previous errors (Wang et al., 2023b). (3) **Paired correct and incorrect actions organized in tree  
083 structures**, to facilitate preference learning. In total, ULTRAINTERACT contains 86K instructions  
084 and 220K action pairs, where each pair consists of an instruction, a correct response, and an incorrect  
085 one. Conceptually, ULTRAINTERACT’s data resemble imbalanced binary trees as shown in Figure 2.

086 ULTRAINTERACT can be used in both supervised fine-tuning and preference learning. Our experi-  
087 ments show that, using ULTRAINTERACT along with established datasets in instruction fine-tuning  
088 already achieves strong performance. ULTRAINTERACT further facilitates preference learning for  
089 reasoning tasks, improving the performance even further with KTO (Ethayarajh et al., 2024) and NCA  
090 (Chen et al., 2024a). Surprisingly, applied to an instruction finetuned EURUS model, DPO (Rafailov  
091 et al., 2023) sometimes hurts the performance.

092 Through careful analysis, we provide evidence that the performance in reasoning correlates with  
093 the value of rewards of chosen data—a higher final reward often indicates a better reasoning  
094 capability. Besides, our investigation suggests that DPO may be less suitable for reasoning tasks  
095 than KTO and NCA. Inspired by this fresh finding, we devise a new objective for reward modeling  
096 to augment the Bradley-Terry objective (Bradley & Terry, 1952), explicitly encouraging training  
097 to increase the absolute rewards of chosen solution and decrease those of rejected data. Furthermore,  
098 ULTRAINTERACT leads to our reward model EURUS-RM-7B, which achieves a better correlation  
099 with human annotators than all existing models on AutoJ (Li et al., 2023a) and MT-Bench (Zheng  
100 et al., 2023), including GPT-4 (OpenAI, 2023). EURUS-RM-7B demonstrates especially strong  
101 preference modeling performance on reasoning tasks.

102 In short, we compiled this work by first synthesizing both SFT and preference datasets to improve  
103 the reasoning ability of open-source models (Section §2). We examined the effectiveness of our  
104 datasets by training both policy and reward models (Section §3). We evaluated the performance of  
105 policy models in Section §4, during which we observed a correlation between reward patterns and  
106 benchmark performances. Next, we then evaluated our reward models and validated that our insights  
107 on the reward-performance correlation can be converted into gains in model training (Section §5).  
Finally, we ablate some factors in our dataset construction in Section §6.

## 2 ULTRAINTERACT: TREE-STRUCTURED ALIGNMENT DATA FOR REASONING

Solving complex problems often requires the model’s capability in planning and reasoning, integrating with tools, and interacting with and learning from both the environment and the users. This is reflected in ULTRAINTERACT’s design choices: (1) Its instructions are diverse, challenging, and of a large scale (§2.1); (2) It provides multi-turn trajectories that solve the input instruction through multiple turns of interaction with and learning from the environment and critique. At each turn, it breaks down the problem into smaller ones (§2.2). (3) ULTRAINTERACT includes pairwise data to facilitate preference learning (§2.3).

Conceptually, ULTRAINTERACT collects a preference tree for each instruction, with the instruction being the root and each action a node (Figure 2). A trajectory is a root-to-leaf path consisting of a sequence of actions. In each preference tree, all nodes of correct actions and all trajectories ending with correct actions can be used for SFT. Paired correct and incorrect nodes or trajectories can be used for preference learning.

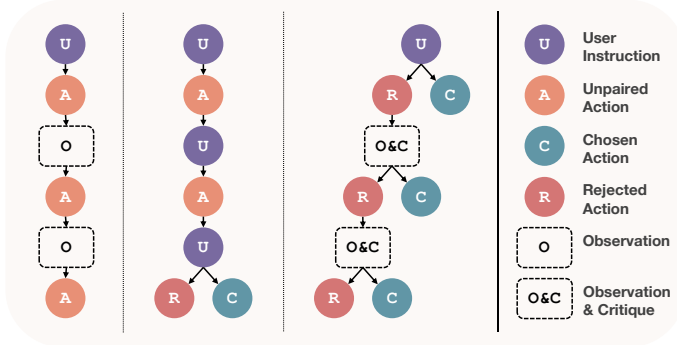


Figure 2: Left: CodeActInstruct (Wang et al., 2024) and Code-Feedback (Zheng et al., 2024); Middle: HH-RLHF (Bai et al., 2022); Right: ULTRAINTERACT. Each instruction in ULTRAINTERACT is constructed as a preference tree.

### 2.1 INSTRUCTION SELECTION EMPHASIZING COMPLEXITY, QUALITY, AND DIVERSITY

We target three representative reasoning tasks: math problem-solving, code generation, and logical reasoning. The complexity, quality, and diversity of the alignment data are crucial to the model’s performance (Liu et al., 2023). Following Wang et al. (2023b), we select challenging problems that GPT-3.5-Turbo fails to solve. We intentionally restrict the selection of the datasets to those with ground-truth solutions, aiming to ensure high-quality oversight signals rather than relying on LLM-as-a-judge annotation (Weyssow et al., 2024). Besides, the gold solutions also serve as references for the critique model to generate feedback. To promote ULTRAINTERACT’s diversity, we pick datasets of different categories. For each dataset, we include distinct reasoning patterns based on question categories or formulations necessary to solve the problems. Table 8 summarizes the datasets selected by ULTRAINTERACT. Except for MATH, none of the training datasets is used in our evaluation.

### 2.2 DECOMPOSITION AND INTERACTION AT EACH TURN

Figure 3 provides an illustrative example. In what follows, we connect the actor model with a Python interpreter as the “environment”. We use GPT-3.5 Turbo as the actor model.

Following Wang et al. (2024), the actor model first decomposes the input problem into several sub-problems and then solves each by generating Python code pieces as actions and using the environment to execute them. To promote solution diversity, the actor model randomly samples one reasoning schema in the form of either CoT (Wei et al., 2022) or modularization programming (Qian et al., 2023; Yuan et al., 2023). The actor then generates actions in text or code to solve each sub-problem, with each step being marked by explicit notations.

Multi-turn interactions with the environment are often necessary to solve challenging problems (Wang et al., 2023b). To improve such capabilities of the models, ULTRAINTERACT collects trajectories in which the actor model interacts with the environment and a critique model (a proxy for user) and refines its action based on their feedback.

The environment receives an action from the actor model along with the interaction history, and then the code interpreter returns two kinds of “Observation”: (1) Python execution results, either program outputs or error traceback messages; (2) binary feedback, indicating whether the solution is

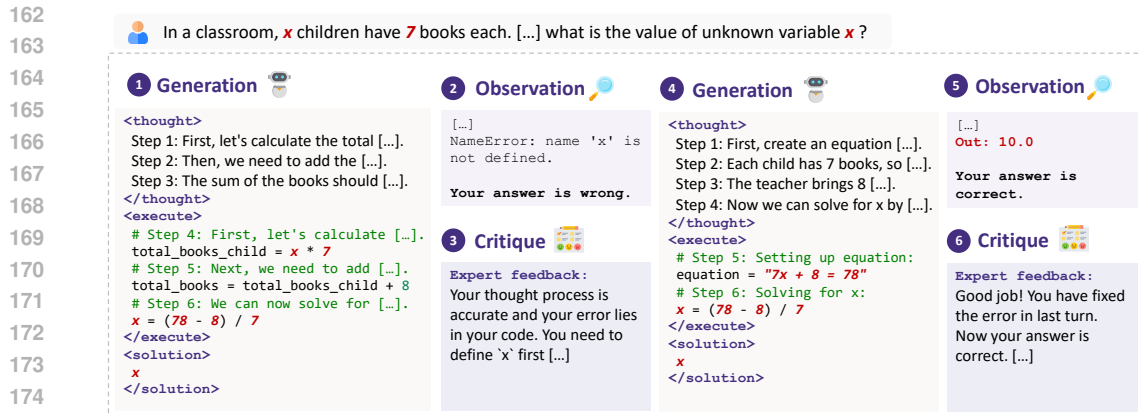


Figure 3: An illustrative example of an ULTRAINTERACT trajectory over two turns. In each turn, the actor model generates step-by-step reasoning chains, and the environment and the critique model provide observations and textual critique respectively.

correct or not. Then, the observations along with the history will be passed to a *critique* model, which locates the errors and provides suggestions for improvements. To avoid potential bias introduced by self-correction (Wang et al., 2023b; Xu et al., 2024), we adopt a stronger model, GPT-4<sup>1</sup>, as the critique and ensure critique quality by providing GPT-4 with ground truth answers as references.

This procedure resembles Wang et al. (2024). However, we adopt more diverse reasoning patterns to teach LLMs to learn rationales rather than simply memorizing answers (Mitra et al., 2023), and learn to create and use tools (Qian et al., 2023; Yuan et al., 2023; Qin et al., 2023). Besides, we believe that it is important for LLMs to learn from the feedback provided by the critique rather than solely from observations of the environment.

### 2.3 PREFERENCE TREES FACILITATES PREFERENCE LEARNING ACROSS MULTIPLE TURNS

Unlike open-ended conversations, where human preference is ambiguous and challenging to specify, many reasoning tasks have clear and objective preferences for *correct* actions. The preference annotation is therefore an evaluation of the correctness of the solutions *conditioning ground truth ones*, which come with the datasets in ULTRAINTERACT. This eliminates the need for human or LLM-based preference annotation and ensures high data quality. To facilitate preference learning, ULTRAINTERACT pairs correct and incorrect actions in each turn.

**Sampling Paired Correct and Incorrect Actions at Each Turn.** For each instruction in ULTRAINTERACT, we sample, from the actor model, a pair of correct and incorrect actions following §2.2. We follow Cui et al. (2023) to sample the pair from different actor models to ensure response diversity. To prevent models from exploiting shortcuts based on surface features, we exclude instances that fail to pass the Python syntax check.

Certain challenging problems in ULTRAINTERACT pose difficulties in obtaining correct actions, even using strong actors such as GPT-4, with nearly zero *pass@100* accuracies. To improve the pass rates of the actor models while keeping the expense under control, we sequentially take the following steps. (1) Directly sampling 20 actions and randomly keeping a correct one, if any. (2) If no correct action is obtained, we repeat the above process up to three times, progressively switching from more cost-effective models to the strong yet expensive GPT-4 Turbo. (3) For the remaining difficult problems where no correct action is acquired after the previous two steps, we provide the actor with ground-truth rationales and answers, and then apply various techniques to elicit correct actions. The specific information provided and the techniques applied vary depending on the tasks (Appendix A.2).

**Tree-structured Action Pairs Across Multiple Turns.** After each turn, the correct action concludes its trajectory. We expand the *incorrect* action into the next turn, and have the actor interact with the environment and the critique to refine its solution (§2.2). We then repeat the procedures introduced earlier in this section to collect an additional action pair. By expanding the *incorrect*

<sup>1</sup>The majority of data collection and experiments are prior to the release of GPT-4o.

action, ULTRAINTERACT can provide data to help models learn from feedback, and collect multiple action pairs for preference learning across multiple turns.

Conceptually, for every instruction, ULTRAINTERACT constructs a binary preference tree with each action being a node (Figure 2). We cap the tree at a maximum of five turns.

**Additional Instruction-action Pairs for Challenging Problems.** We believe the challenging instructions that make it to step (3) above can provide valuable training signals. Therefore, for a subset of these problems with multiple ground truth solutions, we further sample additional correct actions to cover all ground truths. Accordingly, we further sample incorrect actions to pair with these additional correct actions, so that they can be used in both supervised fine-tuning and preference learning.

With preference trees, ULTRAINTERACT enables comparisons at every turn, in contrast to comparing only at the last turn (Bai et al., 2022), and thus can improve the models’ interaction ability. Closing this section, Table 1 summarizes statistics of ULTRAINTERACT. Check more details in Appendix A.4. We show examples of the data in Appendix G.

Table 1: Some statistics of ULTRAINTERACT.

Task	Type		# Instructions	# Turns per Traj.					# Tokens per Traj.	Avg. # Traj per Ins.	Total # Pairs	# Correct Answers
	w/ Interaction?	w/ Tool?		T1	T2	T3	T4	T5				
Math	✓	✓	22,928	10,440	4,122	1,898	904	5,564	1,750.0	1.0	42,780	68,033
	✗	✓	2,757	16,154	-	-	-	-	439.1	5.9	13,217	16,154
	✓	✗	22,639	10,708	3,521	1,459	723	6,228	1,521.9	1.0	44,750	62,182
	✗	✗	2,083	16,348	-	-	-	-	538.1	7.8	12,624	16,348
Coding	✓	-	20,463	13,265	2,584	987	379	3,248	1,728.5	1.0	18,106	22,215
	✗	-	8,495	92,618	-	-	-	-	1,070.4	5.5	78,634	92,618
Logic	✓	✓	2,086	1,685	298	72	8	23	1,299.8	1.0	1,750	2,198
	✓	✗	4,467	2,453	1,674	340	0	0	1,266.7	1.0	7,958	7,231
<b>Total</b>	-	-	85,918	163,671	12,199	4,756	2,014	15,063	1,201.8	2.3	219,819	286,979

### 3 EURUS: STATE-OF-THE-ART OPEN LLMs IN REASONING

ULTRAINTERACT helps us develop EURUS, a suite of LLMs and a reward model (RM).

**Supervised Fine-Tuning.** EURUS-7B-SFT and LLAMA-3-EURUS-8B-SFT are fine-tuned from Mistral-7B (Jiang et al., 2023a) and Llama-3-8B (Meta, 2024) respectively, and EURUS-8x22B-SFT from Mixtral-8x22B (Jiang et al., 2024). First, we perform SFT using all correct actions (287K) in ULTRAINTERACT. We find it yields better performance to discard interaction history and train only on correct leaf nodes in each tree. To improve general instruction-following ability, we include some open-source SFT datasets in our data mixture. Please find mixture details in Appendix B.

**Preference Learning.** Based on EURUS-SFT models, we explore three preference learning algorithms, DPO (Rafailov et al., 2023), KTO (Ethayarajh et al., 2024), and NCA (Chen et al., 2024a). Differently from SFT, here we include all multi-turn trajectory pairs in our ULTRAINTERACT (220K) and include all UltraFeedback (Cui et al., 2023) pairs (340K).

**Reward Modeling.** Similarly to the preference learning, we use all 220K multi-turn trajectory pairs from ULTRAINTERACT; it is further augmented with the 240K single-turn action pairs from ULTRAINTERACT. More details are in the Appendix B. We include all 340K pairs from UltraFeedback and one pair for each instruction from UltraSafety (Guo et al., 2024b), totaling 3K. EURUS-RM-7B is initialized from EURUS-7B-SFT with a new linear layer.

Our findings in §4.2 indicate that the absolute values of rewards make a big difference in the models’ reasoning performance, with decreasing rewards of chosen actions possibly resulting in a performance degradation. We therefore augment the established Bradley-Terry (BT) objective  $\mathcal{L}_{BT}$  with an additional term  $\mathcal{L}_{DR}$  to directly increase the reward of the chosen actions for instances from ULTRAINTERACT, and decrease those of the rejected ones:

$$\mathcal{L}_{ULTRAINTERACT} = \underbrace{-\log\left(\sigma\left(r_{\theta}(x, y_c) - r_{\theta}(x, y_r)\right)\right)}_{\mathcal{L}_{BT}: \text{optimize relative rewards}} - \underbrace{\log\left(\sigma\left(r_{\theta}(x, y_c)\right)\right) - \log\left(\sigma\left(-r_{\theta}(x, y_r)\right)\right)}_{\mathcal{L}_{DR}: \text{increase } r_{\theta}(x, y_c) \text{ and decrease } r_{\theta}(x, y_r)}$$

We train ULTRAINTERACT examples with  $\mathcal{L}_{ULTRAINTERACT}$ , while for instances from other datasets, we train with  $\mathcal{L}_{BT}$  only.  $\theta$  denotes the reward model’s parameters,  $r_{\theta}(x, y_c)$  and  $r_{\theta}(x, y_r)$  the rewards on the chosen and rejected actions respectively. We ablate the importance of  $\mathcal{L}_{BT}$  and  $\mathcal{L}_{DR}$  in §6.2.

Table 2: Open-source LLM baselines that we compare to.

Type	Models
General Purpose	Mistral-7B-Instruct-v0.2 (Jiang et al., 2023a), Zephyr-7B- $\beta$ (Tunstall et al., 2023), OpenChat-3.5-1210 (Wang et al., 2023a), Starling-LM-7B- $\alpha$ (Zhu et al., 2023), Mixtral-8x7B-Instruct (Jiang et al., 2023a), DeepSeek-LLM-67B-Chat (DeepSeek-AI, 2024), QWen1.5-72B-Chat (Bai et al., 2023)
Coding	MagiCoder-S-DS-6.7B (Wei et al., 2023), OpenCodeInterpreter (OpenCI for short, DS-6.7B/CL-70B) (Zheng et al., 2024), DeepSeek-Coder-33B-Instruct (Guo et al., 2024a), and CodeLLaMA-70B-Instruct (Roziere et al., 2023).
Math	MAmmoTH-7B-Mistral (Yue et al., 2023), WizardMath-7B-v1.1 (Luo et al., 2023a), OpenMath (Mistral-7B/CodeLLaMA-70B) (Toshniwal et al., 2024).

## 4 EVALUATION OF EURUS MODELS

**Evaluation Setup.** We consider both single-turn and multi-turn reasoning. For single-turn evaluation, we consider HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and LeetCode (Guo et al., 2024a) for coding, GSM-Plus (Li et al., 2024), MATH, TheoremQA (Chen et al., 2023), SVAMP (Patel et al., 2021), and ASDiv (Miao et al., 2020) for math, and BBH-Hard (Suzgun et al., 2022) for reasoning. We evaluate with pass@1 accuracy. We also use IFEval (Zhou et al., 2023) to assess the instruction-following ability and report the prompt-level loose score. For multi-turn evaluation, we adopt MINT (Wang et al., 2023b) and only consider the coding and math problems. We report the success rate at Turn 5. Please find further details on evaluation setups and evaluations beyond reasoning in Appendix C. As shown in Table 2, we compare our EURUS with general-purpose models, and those specialized in coding and math of various sizes. We also summarize the results of GPT-3.5 Turbo and GPT-4 reported in previous works.

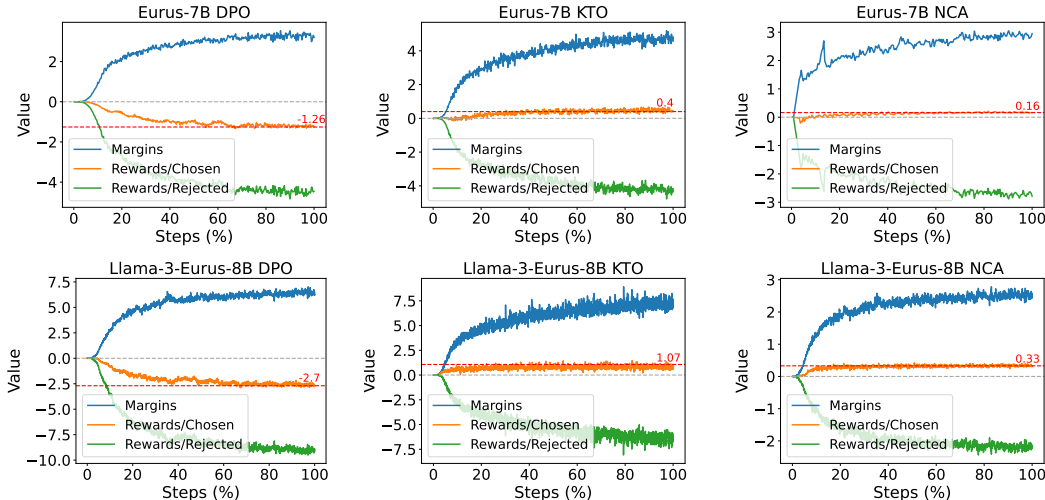
Table 3: Overall performance. All test sets except MATH are out-of-distribution to our models and most baselines. MAmmoTH, OpenChat, and Starling-LM have been trained on TheoremQA test sets. We ~~strikerthrough~~ the contaminated numbers.

Model	Coding			Math				Reasoning	Ins-Following	Multi-Turn		Avg.	
	HumanE.	MBPP	LeetC.	GSM-Plus	MATH	Theo.QA	SVAMP	ASDiv	BBH	IFEval	Code		Math
	~7B												
Mistral-7B-Instruct-v0.2	39.0	30.8	6.1	15.7	9.5	8.5	42.9	49.5	62.4	44.4	7.4	26.2	28.5
Zephyr-7B- $\beta$	29.3	35.8	2.2	23.3	5.0	7.8	19.1	28.0	61.8	39.7	5.2	16.9	22.8
OpenChat-3.5-1210	64.0	61.7	11.7	46.7	28.1	49.4	75.4	77.0	67.0	50.3	21.3	32.4	46.2
Starling-LM-7B- $\alpha$	46.3	51.1	8.9	23.7	21.5	42.0	26.3	39.8	67.1	26.1	18.4	28.9	30.8
MagiCoder-S-DS-6.7B	75.6	<b>70.4</b>	<b>23.9</b>	16.4	19.9	13.1	61.6	62.8	57.0	21.1	<b>27.9</b>	8.0	38.1
OpenCI-DS-6.7B	<b>76.8</b>	66.2	16.1	41.5	31.6	16.1	74.5	79.8	53.9	22.6	5.9	1.3	40.5
MAmmoTH-7B-Mistral	24.4	42.4	7.2	40.1	36.0	<del>26.3</del>	60.7	72.3	57.7	34.9	3.7	6.7	<del>34.4</del>
WizardMath-7B-v1.1	50.0	53.9	6.7	54.6	30.0	16.5	57.8	73.5	64.4	22.6	16.2	8.9	37.9
OpenMath-Mistral-7B	33.5	46.6	11.7	<b>59.4</b>	<b>39.1</b>	13.1	83.4	79.8	58.6	15.0	2.9	5.3	37.4
EURUS-7B-SFT	55.5	59.1	20.0	52.1	32.6	20.0	82.2	84.1	64.6	44.0	15.4	28.4	46.5
+ DPO	50.6	52.1	8.3	51.0	28.3	20.9	78.7	83.8	65.0	42.5	20.6	32.4	44.5
+ KTO	56.1	58.6	18.9	55.0	33.2	20.6	84.4	85.0	67.6	43.1	19.1	<b>43.6</b>	48.8
+ NCA	55.5	60.2	14.4	54.9	34.2	20.9	<b>84.6</b>	85.4	64.3	42.7	21.3	38.7	48.1
LLAMA-3-EURUS-8B-SFT	51.2	57.9	17.2	50.7	32.0	21.3	82.2	83.7	<b>72.4</b>	47.1	18.4	24.5	46.6
+ DPO	43.9	50.1	11.7	45.3	26.8	21.4	54.1	67.5	71.3	<b>56.7</b>	21.3	39.2	42.4
+ KTO	51.8	58.1	15.6	54.8	34.2	24.9	80.1	86.7	71.7	50.6	26.5	37.4	49.4
+ NCA	50.6	60.4	15.6	55.2	34.8	<b>25.4</b>	79.9	<b>87.5</b>	71.7	56.2	21.3	36.3	<b>49.6</b>
	~40B												
Mixtral-8x7B-Instruct	50.6	50.1	5.6	<b>49.6</b>	<b>25.9</b>	20.4	66.4	68.8	<b>73.5</b>	<b>48.8</b>	12.5	<b>37.3</b>	42.5
DeepSeek-Coder-33B-Ins	<b>82.3</b>	<b>73.9</b>	<b>27.8</b>	29.5	20.2	<b>21.9</b>	<b>75.2</b>	<b>85.0</b>	61.5	26.1	<b>35.3</b>	21.8	<b>46.7</b>
	~70B												
CodeLLaMA-70B-Instruct	56.7	58.6	14.4	34.9	12.0	8.4	63.5	70.1	74.5	24.0	3.7	14.2	36.3
DeepSeek-LM-67B-Chat	70.7	65.7	20.0	65.0	41.0	17.9	74.0	84.0	78.9	52.7	30.9	41.8	53.5
QWen1.5-72B-Chat	71.3	56.9	15.6	65.4	43.4	18.5	79.5	79.1	78.0	53.4	27.2	38.2	52.2
OpenCI-CL-70B	<b>77.4</b>	<b>71.7</b>	20.0	46.1	29.2	18.8	76.1	79.4	66.7	26.8	30.9	12.0	46.3
OpenMath-CL-70B	39.0	52.6	15.0	62.2	45.9	15.9	86.6	82.8	59.9	15.7	14.0	0.4	40.8
WizardLM-2-8x22B	72.0	64.2	24.4	57.0	<b>50.9</b>	32.0	81.2	82.2	85.3	<b>68.9</b>	13.2	43.2	56.2
Mixtral-8x22B-Instruct-v0.1	76.2	68.7	23.3	51.2	49.6	33.1	88.3	<b>92.4</b>	<b>86.0</b>	67.1	<b>39.0</b>	59.7	61.2
EURUX-8x22B-KTO	71.3	68.9	29.4	<b>68.3</b>	48.1	35.3	91.5	90.6	83.6	67.1	38.2	57.5	62.5
EURUX-8x22B-NCA	75.0	69.7	<b>35.0</b>	68.1	49.0	<b>35.5</b>	<b>91.6</b>	92.1	83.5	67.1	33.1	<b>63.0</b>	<b>63.6</b>
	Proprietary Models												
GPT-3.5 Turbo	76.8	82.5	23.3	61.2	37.8	35.6	83.0	90.6	70.1	56.6	29.4	36.9	57.0
GPT-4	85.4	<b>83.5</b>	<b>41.8</b>	<b>85.6</b>	<b>69.7</b>	<b>52.4</b>	<b>94.8</b>	<b>92.6</b>	<b>86.7</b>	<b>79.7</b>	<b>59.6</b>	<b>65.8</b>	<b>74.8</b>

### 4.1 RESULTS

According to Table 3, all EURUS variants achieve the best overall performance among open-source models of similar sizes. EURUS even outperforms specialized models in corresponding domains in many cases. Notably, EURUS-7B and LLAMA-3-EURUS-8B outperform baselines that are 5 $\times$  larger, and EURUX-8x22B achieves better performance than GPT-3.5 Turbo. EURUS’s instruction-following performance is among the best general-purpose models, substantially better than specialized ones.

324 **Preference learning with ULTRAINTERACT can further improve the performance, especially in**  
 325 **math and the multi-turn ability.** KTO and NCA consistently improve the models’ performance in  
 326 all five math benchmarks and multi-turn evaluations, while their effects vary in others. Since SFT  
 327 models only use the single-turn data from ULTRAINTERACT while preference learning uses the  
 328 multi-turn ones, the improvements in interaction ability should also be attributed to ULTRAINTERACT  
 329 rather than the algorithms alone (See Appendix H for data mixture ablation results). Surprisingly, we  
 330 observe that **DPO is not as effective as its variants**. We analyze this phenomenon in §4.2.



337 Figure 4: Reward patterns of EURUS-7B and LLAMA-3-EURUS-8B preference learning with DPO,  
 338 KTO, and NCA. For all algorithms, the rewards of rejected data keep decreasing and the margins  
 339 between chosen and rejected data keep increasing. However, the rewards of chosen data decrease  
 340 below zero in DPO while keeping increasing and staying positive in KTO and NCA. The absolute  
 341 values of the reward in the last step (in red) of the three algorithms positively correlate with their  
 342 performance in Table 3.  
 343  
 344  
 345  
 346

347 **4.2 EXPLICIT REWARD AS A PROXY? HYPOTHESIS FOR PREFERENCE LEARNING IN**  
 348 **REASONING**

349 We investigate the reason why DPO behaves differently than KTO and NCA. We start by empirically  
 350 inspecting the rewards throughout the preference learning process, as shown in Figure 4. Rewards for  
 351 chosen and rejected data both keep decreasing through DPO, though the rewards for chosen data is  
 352 still higher hence the loss decreases. In KTO and NCA, the rewards of chosen data keep increasing  
 353 with those of rejected data decreasing.

354 Therefore, we hypothesize it is the distinction in the trend of rewards that leads to the performance gap  
 355 between DPO and the other two algorithms. This distinction can be attributed to that DPO, derived  
 356 from the Bradley-Terry model, only optimizes the relative differences between chosen and rejected  
 357 data overlooking the absolute values of the rewards. This is a non-issue in alignment with general hu-  
 358 man values where preference is “relative” and there can be many valid answers to the same input. How-  
 359 ever, in reasoning tasks, the space of correct answers is much smaller than that of incorrect ones. Fur-  
 360 ther, we notice that the rewards of chosen data in the last training step follow the ranking order of KTO  
 361 > NCA > DPO, positively correlate with their performance trends. Therefore, we believe that increas-  
 362 ing the rewards of the chosen data is especially beneficial in preference learning for reasoning tasks.  
 363  
 364  
 365  
 366  
 367  
 368  
 369

370 **5 EVALUATION OF EURUS-RM-7B**

371 **5.1 REWARD MODELING PERFORMANCE**

372 We evaluate EURUS-RM-7B on three RM benchmarks, RewardBench (Lambert et al., 2024), AutoJ  
 373 (Li et al., 2023a), and MT-Bench (Zheng et al., 2023). Aiming for a more realistic OOD evaluation,  
 374 we exclude the “prior sets” split from RewardBench, since many baselines train on the datasets that  
 375 this split contains. We compare with PairRM (Jiang et al., 2023b), Starling-RM-7B/34B (Zhu et al.,  
 376 2023), UltraRM-13B (Cui et al., 2023), GPT-3.5 Turbo, and GPT-4.  
 377

Table 4: Results on reward modeling benchmarks. UF: UltraFeedback; US: UltraSafety. The best performance in each benchmark is in bold and the second best one is underlined. Most baseline results are from Jiang et al. (2023b) and Lambert et al. (2024).

Model	Reward Bench					AutoJ				MT-Bench
	Chat	Chat-Hard	Safety	Reasoning	Avg.	Code	Math	Others	Overall	
PairRM	90.2	53.0	31.5	60.0	58.7	58.3	52.8	58.9	59.1	59.0
Starling-RM-7B	98.0	43.4	88.6	74.6	76.2	59.2	47.2	61.4	60.8	56.8
Starling-RM-34B	96.9	59.0	89.9	90.3	<b>84.0</b>	65.8	54.2	62.3	62.6	60.4
UltraRM-13B	96.1	55.3	45.8	82.0	69.8	55.0	43.1	59.6	59.9	56.0
GPT-3.5 Turbo	-	-	-	-	-	36.6	40.3	41.2	42.7	57.1
GPT-4	-	-	-	-	-	69.2	51.4	61.4	61.9	63.9
EURUS-RM-7B	96.5	65.3	80.7	87.0	<u>82.4</u>	87.5	82.5	78.0	80.7	<u>79.4</u>
w/o $\mathcal{L}_{DR}$	96.4	59.9	79.5	77.5	78.3	83.8	82.5	78.9	80.7	79.3
w/o $\mathcal{L}_{BT}$	96.8	58.5	83.8	84.2	80.8	88.8	92.5	79.4	<b>81.9</b>	<b>79.6</b>
w/o US	96.5	66.2	67.7	81.7	73.3	87.5	90.0	79.2	<u>81.8</u>	79.2
w/o UF + US	95.1	61.1	63.7	73.4	78.0	73.8	80.0	71.7	72.8	73.0

**Results.** Table 4 summarizes reward modeling performance. Takeaways are as follows:

**EURUS-RM-7B stands out as the best 7B RM overall, and achieves similar or better performance than much larger baselines. Particularly, it outperforms GPT-4 in certain tasks.** EURUS-RM-7B achieves a better correlation with human experts than all existing models on AutoJ and MT-Bench, and it achieves comparable performance to the 5× larger Starling-RM-34B on RewardBench. On RewardBench, EURUS-RM-7B outperforms all baselines on the “Chat-Hard” split while achieving very competitive performance on the “Reasoning” split. Across the AutoJ splits, EURUS-RM-7B outperforms all baselines, with the only exception being GPT-4’s results on Coding.

**Our training objective is beneficial in improving RM performance on hard problems and reasoning.** Table 4 shows that optimizing  $\mathcal{L}_{DR}$  improves RM’s reasoning ability, but BT modeling is still beneficial in equipping RM with abilities in general chatting as suggested in the “Chat-Hard” column, though its effect on reasoning may vary.

**ULTRAINTERACT is compatible with other datasets like UltraFeedback and UltraSafety, and mixing these datasets can balance different RM abilities.** Improving RM’s capabilities in reasoning with ULTRAINTERACT does not sacrifice others, which indicates that ULTRAINTERACT can be a great ingredient for the training data mixture of reward models.

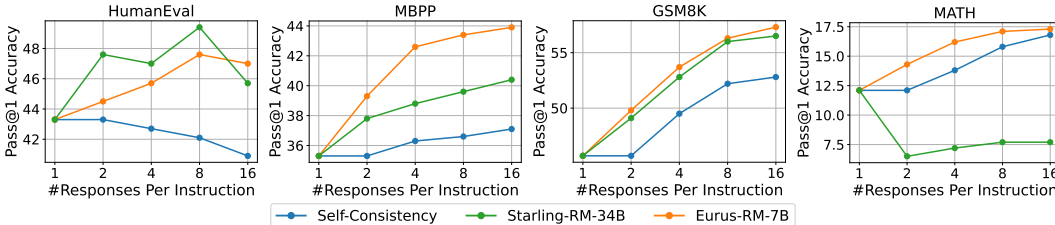


Figure 5: Results on reranking Mistral-7B-Instruct-v0.2’s responses. Full results in Table 11.

## 5.2 BEST-OF-N AND MCTS DECODING

To further explore EURUS-RM-7B’s potential in improving models’ performance through reranking, we use it to rerank Mistral-7B-Instruct-v0.2’s responses on HumanEval, MBPP, GSM8K, and MATH. We report the results of random sampling, self-consistency, and Starling-RM-34B as baselines. Finally, we examine the ability of EURUS-RM-7B to guide Mistral-7B-Instruct-v0.2 for MCTS decoding on math, compared to Starling-RM-7B. Due to the inference overhead, we only sample 500 samples from GSM8K and MATH for evaluation.

**Results.** Figure 5 and Table 5 presents some results with others in Appendix D.1.

**EURUS-RM-7B improves LLMs’ reasoning performance by a large margin through reranking.** It consistently improves pass@1 accuracy across all tasks and performs better than 5× larger baseline Starling-RM-34B. Also, EURUS-RM-7B’s reranking performance scales well with #responses per



instruction, except for a slight decrease in HumanEval when increasing the response number from 8 to 16. In contrast, Starling-RM-34B suffers from a severe performance drop on HumanEval and it consistently hurts model accuracy on MATH.

**EURUS-RM-7B can improve the performance of Mistral-7B-Instruct-v0.2 by 81.4% on GSM8K and 43.8% on MATH.** According to Table 5, the policy model only achieves 25.8% and 6.4% on GSM8K and MATH respectively, but the performances are largely improved with RM-guided MCTS decoding. Compared to Starling-RM-7B, EURUS-RM-7B shows an advantage of 46.8% vs. 44.8% on GSM8K, and 9.2% vs. 6.6% on MATH.

Table 5: Mistral-7B-Instruct-v0.2 results with RM-guided MCTS decoding.

Model	GSM8K	MATH
Mistral-7B-Instruct-v0.2	25.8	6.4
+ Starling-RM-7B	44.8	6.6
+ EURUS-RM-7B	46.6	9.2

### 5.3 HOW DOES $\mathcal{L}_{DR}$ WORK IN REWARD MODELING?

Results in Table 4 have demonstrated the effectiveness of our proposed reward modeling objective. To figure out the working mechanism, we plot the reward patterns in reward modeling in Figure 6. As we can see,  $\mathcal{L}_{BT}$  only optimize the reward margin, and the reward trend of chosen data is unstable while rewards of rejected data are positive for the most time. In contrast,  $\mathcal{L}_{DR}$  forces the reward of chosen data to be positive and that of rejected data to be negative as expected. There is a consistent trend of increasing the reward of chosen data and decreasing that of rejected data, which also leads to a widening reward margin. When combining  $\mathcal{L}_{BT}$  with  $\mathcal{L}_{DR}$ , the resulting pattern is more similar to using  $\mathcal{L}_{DR}$  only. These results indicate the importance of the absolute value of rewards in reward model training.

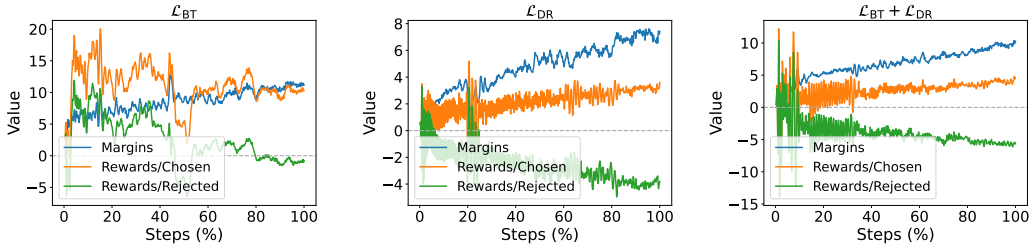


Figure 6: Reward pattern in reward modeling. The rewards of chosen data and margins increase regardless of  $\mathcal{L}_{DR}$ , but the rewards of rejected data decrease to be negative with regularization.

## 6 ANALYSIS

### 6.1 ARE PROPRIETARY MODELS NECESSARY FOR HIGH-QUALITY DATA CONSTRUCTION?

The current construction of ULTRAINTERACT largely adopts GPT models. However, using proprietary models raises concerns on the non-permissive license and the heavy burden in financial budgets, both of the which can be addressed by open-source models. Therefore, to ablate the use of proprietary models and examine how far we can go with open-source models for data synthesis, we supplement another version of ULTRAINTERACT by substituting all GPT responses with those generated by open-source models, **including EURUS itself**. We list the statistics, including models used to sample responses and corresponding proportions, in Appendix F. We named the original GPT-involved version as ULTRAINTERACT-v1 and the open-source model generated as ULTRAINTERACT-v2. We retrain LLAMA-3-EURUS-8B with ULTRAINTERACT-v2 and union of v1 and v2. Results are presented in Table 6.

Surprisingly, LLAMA-3-EURUS-8B performances are greatly improved. Compared to v1, training on v2 improves model performance on both SFT and preference learning stage, and particularly, LLAMA-3-EURUS-8B-KTO (v2) successfully surpasses the official Llama-3-8b-Instruct model (Meta, 2024), which it previously failed to, indicating we can train strong models without distilling proprietary models. Further, trained on the mix of v1 and v2, LLAMA-3-EURUS-8B consistently outperforms Llama-3-8B-Instruct after preference learning. Notably, to the best of our knowledge, this is the first open recipe on Llama-3-8B that can outperform the official Llama-3-8B-Instruct model on reasoning.

### 6.2 ABLATION STUDY

Table 6: Performances of LLAMA-3-EURUS-8B trained with different versions of ULTRAINTERACT.

Model	Coding			Math					Reasoning	Ins-Following	Multi-Turn		Avg.
	HumanEval	MBPP	LeetC.	GSM-Plus	MATH	Theo.QA	SVAMP	ASDiv	BBH	IFEval	Code	Math	
LLAMA-3-8b-Instruct	59.8	60.4	21.1	49.1	28.3	15.0	73.7	75.9	73.6	81.9	25.0	44.7	50.7
LLAMA-3-EURUS-8B-SFT (v1)	51.2	57.9	17.2	50.7	32.0	21.3	82.2	83.7	72.4	47.1	18.4	24.5	46.6
+ DPO	43.9	50.1	11.7	45.3	26.8	21.4	54.1	67.5	71.3	56.7	21.3	39.2	42.4
+ KTO	51.8	58.1	15.6	54.8	34.2	24.9	80.1	86.7	71.7	50.6	26.5	37.4	49.4
+ NCA	50.6	60.4	15.6	55.2	34.8	25.4	79.9	87.5	71.7	56.2	21.3	36.3	49.6
LLAMA-3-EURUS-8B-SFT (v2)	55.5	60.2	17.8	54.4	37.7	24.6	88.1	85.0	73.2	49.9	18.4	29.7	49.5
+ DPO	57.9	53.9	12.2	49.1	37.5	27.6	78.6	76.9	73.1	55.6	22.8	44.0	49.1
+ KTO	59.1	55.9	21.1	60.1	39.7	25.9	86.9	86.6	73.4	47.5	22.1	38.8	51.4
+ NCA	56.7	52.6	15.6	56.8	36.2	24.1	88.0	84.1	73.6	51.0	21.3	36.3	49.7
LLAMA-3-EURUS-8B-SFT (v1 + v2)	54.9	62.2	16.7	54.4	38.0	25.0	86.4	85.9	72.5	59.1	21.3	30.4	50.6
+ DPO	53.0	54.4	18.3	58.8	37.0	27.5	83.1	89.0	72.2	51.0	25.0	42.1	51.0
+ KTO	61.6	59.1	15.6	61.2	39.0	26.8	86.8	88.2	72.1	47.7	29.4	39.2	52.2
+ NCA	53.7	60.9	12.8	62.2	38.1	25.4	87.3	87.7	72.5	49.0	25.0	35.5	50.8

We study the impact of ULTRAINTERACT and other open-source alignment data on EURUS-7B-SFT’s performance. We consider three settings: (1) **With original ground-truth answers**, which replaces the generated actions with ground-truth rationales and answers from the original datasets. If no rationales are available, we use those from ULTRAINTERACT. (2) **Existing data only**. (3) **ULTRAINTERACT only**. We evaluate with the same setting as §4 and report the averaged scores. See full results in Appendix E.

Table 7: Ablation study of SFT data.

Model	Coding	Math	BBH	IFEval	Avg.
EURUS-7B-SFT	44.9	58.5	64.6	44.0	53.6
Ground-truth	33.9	46.1	64.4	42.9	44.0
Existing Data Only	31.2	33.5	65.3	43.6	37.0
ULTRAINTERACT Only	37.3	56.2	67.0	17.4	47.7

In Table 7, EURUS outperforms the “Ground-truth” model on all tasks, confirming the advantage of ULTRAINTERACT’s designs of divide-and-conquer and code-as-action patterns, in line with conclusions of concurrent work (Chen et al., 2024b; Wang et al., 2024). Training only on existing data without ULTRAINTERACT greatly hurts the reasoning performance, confirming the effectiveness of ULTRAINTERACT. Meanwhile, training only on ULTRAINTERACT suffers a performance drop except for BBH, especially in instruction following. We attribute the performance drop to a worse instruction-following ability. This suggests the necessity of mixing ULTRAINTERACT with other alignment data for better all-around supervised fine-tuning.

## 7 RELATED WORK

**Open LLMs in Reasoning.** Open-source LLMs have shown remarkable progress in building *specialists* that excel in mathematics reasoning (Luo et al., 2023a; Yue et al., 2023; Toshniwal et al., 2024) or coding abilities (Roziere et al., 2023; Wei et al., 2023; Guo et al., 2024a; Zheng et al., 2024). On the contrary, mastering general reasoning capabilities still challenges open models, while the most advanced ones (DeepSeek-AI, 2024; Bai et al., 2023; Touvron et al., 2023; Jiang et al., 2024) are well behind proprietary models. More, these cutting-edge open general-purpose models maintain their alignment recipes confidential, which further hinders the replication and development of open-source reasoning models.

**Preference Learning for Reasoning.** Aligning language models from human or AI preferences has emerged as a prevalent approach in the open-source community (Tunstall et al., 2023; Bai et al., 2023) with the proposal of DPO (Rafailov et al., 2023) and high-quality preference datasets (Cui et al., 2023; Zhu et al., 2023). Different from open-domain chatbots, preference learning is largely underexplored in complex reasoning. Recent research showed performance degradation when applying DPO on reasoning tasks, but some newly proposed algorithms demonstrated a positive effect (Ethayarajh et al., 2024; Chen et al., 2024a; Mitra et al., 2024; Shao et al., 2024b). However, a deep understanding of preference learning, specifically its efficacy on complex reasoning, is not yet established.

## 8 CONCLUSION

We strive to narrow the huge gap between open-source models and proprietary models from the perspective of alignment. Our work pushes the boundaries of open-source reasoning generalists by (1) releasing a high-quality multi-turn reasoning dataset ULTRAINTERACT with preference trees, (2) introducing EURUS-series LLMs which achieve new SOTA on challenging reasoning benchmarks and (3) providing insights on preference learning for reasoning through analysis, leading to new reward modeling objectives as well as a powerful reward model for reasoning.

## REFERENCES

- 540  
541  
542 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh  
543 Hajishirzi. MathQA: Towards interpretable math word problem solving with operation-based  
544 formalisms. In *Proc. of NAACL-HLT*, 2019.
- 545 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
546 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
547 models. *ArXiv preprint*, abs/2108.07732, 2021.
- 548 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
549 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,  
550 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan,  
551 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin  
552 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng  
553 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren  
554 Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *ArXiv preprint*, abs/2309.16609,  
555 2023.
- 556 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn  
557 Drain, Stanislav Fort, Deep Ganguli, T. J. Henighan, Nicholas Joseph, Saurav Kadavath, John  
558 Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez,  
559 Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario  
560 Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and  
561 Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human  
562 feedback. *ArXiv preprint*, abs/2204.05862, 2022.
- 563 Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method  
564 of paired comparisons. *Biometrika*, 39, 1952.
- 565  
566 Huayu Chen, Guande He, Hang Su, and Jun Zhu. Noise contrastive alignment of language models  
567 with explicit rewards. *ArXiv preprint*, abs/2402.05369, 2024a.
- 568 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
569 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
570 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
571 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,  
572 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios  
573 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino,  
574 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,  
575 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,  
576 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob  
577 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
578 large language models trained on code, 2021.
- 579 Wenhu Chen, Ming Yin, Max W.F. Ku, Yixin Wan, Xueguang Ma, Jianyu Xu, Tony Xia, Xinyi  
580 Wang, and Pan Lu. Theoremqa: A theorem-driven question answering dataset. *ArXiv preprint*,  
581 abs/2305.12524, 2023.
- 582 Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and  
583 Feng Zhao. Agent-flan: Designing data and methods of effective agent tuning for large language  
584 models. volume abs/2403.12881, 2024b.
- 585  
586 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
587 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
588 math word problems. volume abs/2110.14168, 2021.
- 589 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu,  
590 and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *ArXiv*  
591 *preprint*, abs/2310.01377, 2023.
- 592  
593 DeepSeek-AI. Deepseek llm: Scaling open-source language models with longtermism. *ArXiv*  
*preprint*, abs/2401.02954, 2024.

- 594 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong  
595 Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional  
596 conversations. In Conference on Empirical Methods in Natural Language Processing, 2023.
- 597  
598 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model  
599 alignment as prospect theoretic optimization. ArXiv preprint, abs/2402.01306, 2024.
- 600  
601 Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. Specializing smaller language  
602 models towards multi-step reasoning. In Proceedings of the International Conference on Machine  
603 Learning, 2023.
- 604  
605 Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle  
606 use a laptop? a question answering benchmark with implicit reasoning strategies. Transactions of  
607 the Association for Computational Linguistics, 9, 2021.
- 608  
609 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
610 Bi, Yu Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When  
611 the large language model meets programming - the rise of code intelligence. ArXiv preprint,  
612 abs/2401.14196, 2024a.
- 613  
614 Yiju Guo, Ganqu Cui, Lifan Yuan, Ning Ding, Jiexin Wang, Huimin Chen, Bowen Sun, Ruobing  
615 Xie, Jie Zhou, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Controllable preference optimization:  
616 Toward controllable multi-objective alignment. ArXiv preprint, abs/2402.19085, 2024b.
- 617  
618 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin  
619 Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence  
620 with apps. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and  
621 Benchmarks Track (Round 2), 2021a.
- 622  
623 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn  
624 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset.  
625 In Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks  
626 Track (Round 2), 2021b.
- 627  
628 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
629 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
630 Mistral 7b. ArXiv preprint, abs/2310.06825, 2023a.
- 631  
632 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
633 Bamford, Devendra Singh Chaplot, Diego de Las Casas, Emma Bou Hanna, Florian Bressand,  
634 Gianna Lengyel, Guillaume Bour, Guillaume Lample, L’elio Renard Lavaud, Lucile Saulnier,  
635 Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak,  
636 Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El  
637 Sayed. Mixtral of experts. 2024.
- 638  
639 Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with  
640 pairwise ranking and generative fusion. In Annual Meeting of the Association for Computational  
641 Linguistics, 2023b.
- 642  
643 Nathan Lambert, Valentina Pyatkin, Jacob Daniel Morrison, Lester James Validad Miranda,  
644 Bill Yuchen Lin, Khyathi Raghavi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi,  
645 Noah A. Smith, and Hanna Hajishirzi. Rewardbench: Evaluating reward models for language  
646 modeling. 2024.
- 647  
648 Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, Hai Zhao, and Pengfei Liu. Generative judge  
649 for evaluating alignment. ArXiv preprint, abs/2310.05470, 2023a.
- 650  
651 Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive  
652 benchmark for evaluating the robustness of llms as mathematical problem solvers. ArXiv preprint,  
653 abs/2402.19255, 2024.
- 654  
655 Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and  
656 Ge Li. Taco: Topics in algorithmic code generation dataset. volume abs/2312.14852, 2023b.

- 648 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom  
649 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien  
650 de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal,  
651 Alexey Cherepanov, James Molloy, Daniel Mankowitz, Esme Sutherland Robson, Pushmeet Kohli,  
652 Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with  
653 alphacode. volume abs/2203.07814, 2022.
- 654 Wing Lian, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and ”Teknium”.  
655 Openorca: An open dataset of gpt augmented flan reasoning traces. [https://https://](https://https://huggingface.co/Open-Orca/OpenOrca)  
656 [huggingface.co/Open-Orca/OpenOrca](https://https://huggingface.co/Open-Orca/OpenOrca), 2023.
- 657 Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for  
658 alignment? a comprehensive study of automatic data selection in instruction tuning. 2023.
- 659 Zihan Liu, Wei Ping, Rajarshi Roy, Peng Xu, Chankyu Lee, Mohammad Shoeybi, and Bryan  
660 Catanzaro. Chatqa: Surpassing gpt-4 on conversational qa and rag. 2024. URL [https://api.](https://api.semanticscholar.org/CorpusID:267035133)  
661 [semanticscholar.org/CorpusID:267035133](https://api.semanticscholar.org/CorpusID:267035133).
- 662 Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter  
663 Clark, and Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured  
664 mathematical reasoning. In Proceedings of ICLR, 2023.
- 665 Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng,  
666 Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical  
667 reasoning for large language models via reinforced evol-instruct. ArXiv preprint, abs/2308.09583,  
668 2023a.
- 669 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
670 Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with  
671 evol-instruct, 2023b.
- 672 Meta. The llama 3 herd of models. ArXiv, 2024.
- 673 Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing  
674 English math word problem solvers. In Proc. of ACL, 2020.
- 675 Swaroop Mishra, Arindam Mitra, Neeraj Varshney, Bhavdeep Sachdeva, Peter Clark, Chitta Baral,  
676 and Ashwin Kalyan. NumGLUE: A suite of fundamental yet challenging mathematical reasoning  
677 tasks. In Proc. of ACL, 2022.
- 678 Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Coda, Clarisse Simoes, Sahaj Agrawal,  
679 Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, Hamid Palangi, Guoqing Zheng,  
680 Corby Rosset, Hamed Khanpour, and Ahmed Awadallah. Orca 2: Teaching small language models  
681 how to reason. ArXiv preprint, abs/2311.11045, 2023.
- 682 Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the  
683 potential of slms in grade school math. ArXiv preprint, abs/2402.14830, 2024.
- 684 OpenAI. Gpt-4 technical report, 2023.
- 685 Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In  
686 Proc. of ACL, 2015.
- 687 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple  
688 math word problems? In Proceedings of the 2021 Conference of the North American Chapter of  
689 the Association for Computational Linguistics: Human Language Technologies, 2021.
- 690 Cheng Qian, Chi Han, Yi Ren Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. Creator: Tool creation  
691 for disentangling abstract and concrete reasoning of large language models. In Conference on  
692 Empirical Methods in Natural Language Processing, 2023.
- 693 Yujia Qin, Shi Liang, Yining Ye, Kunlun Zhu, Lan Yan, Ya-Ting Lu, Yankai Lin, Xin Cong, Xiangru  
694 Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Marc H. Gerstein, Dahai Li,  
695 Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+  
696 real-world apis. ArXiv preprint, abs/2307.16789, 2023.

- 702 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea  
703 Finn. Direct preference optimization: Your language model is secretly a reward model. ArXiv  
704 preprint, abs/2305.18290, 2023.
- 705 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
706 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code.  
707 ArXiv preprint, abs/2308.12950, 2023.
- 708 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu,  
709 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language  
710 models, 2024a. URL <https://arxiv.org/abs/2402.03300>.
- 711 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Y Wu,  
712 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language  
713 models. ArXiv preprint, abs/2402.03300, 2024b.
- 714 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
715 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. Challenging  
716 big-bench tasks and whether chain-of-thought can solve them. ArXiv preprint, abs/2210.09261,  
717 2022.
- 718 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Git-  
719 man. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. arXiv preprint arXiv:  
720 Arxiv-2402.10176, 2024.
- 721 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
722 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas  
723 Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,  
724 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S.  
725 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,  
726 Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril,  
727 Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar  
728 Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan  
729 Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross  
730 Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang,  
731 Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey  
732 Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. ArXiv  
733 preprint, abs/2307.09288, 2023.
- 734 Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,  
735 Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct  
736 distillation of lm alignment. ArXiv preprint, abs/2310.16944, 2023.
- 737 Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. Openchat:  
738 Advancing open-source language models with mixed-quality data. ArXiv preprint, abs/2309.11235,  
739 2023a.
- 740 Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. Mint:  
741 Evaluating llms in multi-turn interaction with tools and language feedback. ArXiv preprint,  
742 abs/2309.10691, 2023b.
- 743 Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji.  
744 Executable code actions elicit better llm agents. ArXiv preprint, abs/2402.01030, 2024.
- 745 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le,  
746 and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. ArXiv  
747 preprint, abs/2201.11903, 2022.
- 748 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is  
749 all you need, 2023.
- 750 Martin Weyssow, Aton Kamanda, and Houari Sahraoui. Codeultrafeedback: An llm-as-a-judge  
751 dataset for aligning large language models to coding preferences. ArXiv preprint, abs/2403.09032,  
752 2024.

756 Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, and William Yang Wang. Perils  
757 of self-feedback: Self-bias amplifies in large language models. [ArXiv preprint, abs/2402.11436](#),  
758 2024.

759 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov,  
760 and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question  
761 answering. In [Proc. of EMNLP](#), 2018.

762

763 Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset  
764 requiring logical reasoning. In [Proc. of ICLR](#), 2020.

765

766 Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi Ren Fung, Hao Peng, and Heng Ji. Craft: Customizing  
767 llms by creating and retrieving from specialized toolsets. [ArXiv preprint, abs/2309.17428](#), 2023.

768 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.  
769 Mammoth: Building math generalist models through hybrid instruction tuning. [ArXiv preprint](#),  
770 [abs/2309.05653](#), 2023.

771 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
772 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haoteng Zhang, Joseph Gonzalez, and Ion Stoica.  
773 Judging llm-as-a-judge with mt-bench and chatbot arena. [ArXiv preprint, abs/2306.05685](#), 2023.

774

775 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and  
776 Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement.  
777 [ArXiv preprint, abs/2402.14658](#), 2024.

778 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny  
779 Zhou, and Le Hou. Instruction-following evaluation for large language models. [ArXiv preprint](#),  
780 [abs/2311.07911](#), 2023.

781

782 Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm  
783 helpfulness & harmlessness with rlaiif, 2023.

784 Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li,  
785 Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models  
786 in code intelligence. [arXiv preprint arXiv:2406.11931](#), 2024.

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

Table 8: ULTRAINTERACT covers a diverse set of datasets spanning three tasks.

Task	Datasets
<b>Math</b>	GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021b), MathQA (Amini et al., 2019), NumGlue (Mishra et al., 2022), TabMWP (Lu et al., 2023)
<b>Coding</b>	CodeContest (Li et al., 2022), TACO (Li et al., 2023b), WikiTableQuestions (Pasupat & Liang, 2015), Magicoder-Evol-Instruct (Luo et al., 2023b; Wei et al., 2023)
<b>Logic</b>	ReClor (Yu et al., 2020), HotpotQA (Yang et al., 2018), StrategyQA (Geva et al., 2021)

## A ADDITIONAL DETAILS IN ULTRAINTERACT CONSTRUCTION

### A.1 DATASET DETAILS

**Math.** We adopt **GSM8K** (Cobbe et al., 2021), **MATH** (Hendrycks et al., 2021b), **MathQA** (Amini et al., 2019), and **NumGLUE** (Mishra et al., 2022) for mathematic reasoning, and include **TabMWP** (Lu et al., 2023) for tabular processing. We retain all the instructions for all datasets except MathQA, NumGLUE, and TabMWP. MathQA divides problems into different categories according to the topics and annotates the formula that indicates the pattern needed to solve each problem. We apply stratified sampling to sample at most five problems for each pattern and prioritize the problems that come from the long-tail category. Numglue contains eight different reasoning tasks and we discard Task 5 (Reading Comprehension + Explicit Numerical Reasoning), Task 6 (Reading Comprehension + Implicit Numerical Reasoning), and Task 7 (Quantitative NLI) due to the simplicity Mishra et al. (2022). For TabMWP, we only keep the questions with difficulty levels 4 and 5 since the rest are too easy for current state-of-the-art models.

**Code.** We focus on programming with Python for the simplicity of integration of the interpreter. We use **CodeContest** (Li et al., 2022) and **TACO** (Li et al., 2023b), two competition-level coding datasets collected from various online platforms. We filter out the overlapped questions. Note that part of the questions in TACO only contain ground-truth solutions and do not contain test cases for evaluation, hence we apply GPT-4 to generate 12 test case inputs (4 basic inputs, 4 edge cases, and 4 large numbers) for each question and then execute the ground-truth solution snippets to produce outputs. Given that the two datasets mainly focus on competition problems that may deviate from real-world daily uses, we exclusively adopt **Magicoder-Evol-Instruct** (Luo et al., 2023b; Wei et al., 2023), the only dataset in our selection that does not contain test cases or ground-truth solutions. We employ GPT-4 Turbo to judge the correctness of generated code during interaction, and therefore we do not use this dataset for preference learning since we cannot rigorously construct pairs of correct and incorrect actions limited by the evaluation reliability. We also include **WikiTableQuestions** (Pasupat & Liang, 2015) for table processing with code.

**Logical Reasoning.** we use the multi-hop reasoning datasets **HotpotQA** (Yang et al., 2018) and **StrategyQA** (Geva et al., 2021), and the logical reasoning dataset **ReClor** (Yu et al., 2020). We follow the setting of Wang et al. (2023b) and convert HotpotQA to a generation task, removing the contexts and requiring LLMs to search relevant information using Wikipedia API.

### A.2 DETAILS ON PREFERENCE TREE CONSTRUCTION

**Models Adopted for Incorrect Action Sampling.** We randomly sample one model from Mistral-7B-Instruct-v0.2, DeepSeek-Coder-33B-Instruct, Mixtral-8x7B-Instruct, and DeepSeek-LLM-67B-Chat to generate one incorrect action to pair with each correct one.

#### **Correct Action Generation Based on Ground Truth Annotations.**

We adopt GPT-3.5 Turbo as the generator to generate correct actions based on ground truth considering the instruction-following ability. We provide different access to the ground truth information for different tasks, specifically: (1) For coding, where test cases are black boxes to reference solutions, we provide full access to the solution codes. The actor model will add step marks and corresponding explanations to the ground-truth code to make it easier to understand, or further refine the code for optimization. (2) For tool-free math problems, to avoid the actor model directly copying the answers to pass the correctness checking, we mask the answer numbers in the rationale before providing it to LLMs. This approach can better ensure response quality since it encourages LLMs to generate



Table 9: Stats breakdown

Task	Dataset	w/ Tool?	# Prompts	# Pairs	# Correct Answers.	Avg. Length	Human Annotation	
							Has Answer?	Has Rationale?
Math	GSM8K	✓	4,522	10,277	17,392	1,746.7	✓	✓
		✗	7,257	10,879	15,752	823.3	✓	✓
	MATH	✓	7,474	22,905	34,667	1,189.0	✓	✓
		✗	7,471	25,765	36,005	1,735.0	✓	✓
	MathQA	✓	7,552	15,079	20,328	2,338.5	✓	✓
		✗	7,159	17,743	22,500	1,916.3	✓	✓
	NumGLUE	✓	3,020	3,601	5,717	1,474.6	✓	✗
✗		2,835	2,987	4,273	1,056.1	✓	✗	
TabMWP	✓	3,117	4,135	6,083	842.6	✓	✗	
Coding	CodeContest	-	8,167	44,319	44,666	2,061.7	✓	✓
	TACO	-	9,016	50,877	58,191	2,143.5	✓	✓
	WikiTableQuestions	-	1,401	1,544	1,738	1,794.8	✓	✗
	Magicoder-Evol-Instruct	-	10,374	0	10,238	687.1	✗	✗
Logic	Reclor	✗	4,467	7,958	7,231	1,266.7	✓	✗
	HotpotQA	✓	1,182	1,009	1,230	1,333.2	✓	✗
	StrategyQA	✓	904	741	968	1,256.2	✓	✗

responses with complete reasoning chains with each step clearly marked. (3) For program-enhanced math reasoning, we first translate the textual rationale into code. Then, we either directly provide it to the actor model to generate plans, or ask the actor model to convert the code into modularization programming and then make plans to create tools to solve problems.

### A.3 DATA DECONTAMINATION

We conduct careful decontamination. Firstly, for LeetCode, we apply the Exact Substring Matching Algorithm<sup>2</sup> to compare with each instruction in the ULTRAINTERACT and find no overlaps. For others, we perform 8-gram exact matching to compare ULTRAINTERACT instructions with test sets of the same task. We remove those instructions that overlap 8 grams with any test sample.

### A.4 DETAILED STATISTICS

In total, ULTRAINTERACT has 86K instructions and 220K action pairs. The Total # Pairs does not equal Total # Turns in ULTRAINTERACT, since we fail to generate sufficient correct actions for every incorrect action in multi-turn trajectories mainly due to a lack of sufficient ground truth annotations. The total # pairs may not equal # correct answers, either, because it is also difficult and unnecessary to sample incorrect actions for the correct ones for some simple instructions. We present the specific information for each dataset. In particular, we list information on human annotation in each dataset, which plays an important role in correct action generation (§2.3 and Appendix A.2). All three steps of correct action sampling methods mentioned in §2.3 can be applied to datasets that have rationales, while for datasets only containing answers, only the first two steps are applicable. We do not apply any of the three-step methods to generate correct answers for Magicoder, the only dataset without any human annotation, to construct preference pairs.

## B ADDITIONAL DETAILS ON TRAINING EURUS MODELS

**Supervised Fine-Tuning.** We finetune base models for 1 epoch with a  $2e-5$  learning rate and 0.1 warmup ratio using a cosine scheduler. For EURUS-7B, we mix 32K UltraChat (Ding et al., 2023), 30K ShareGPT<sup>3</sup>, and 50K OpenOrca(Lian et al., 2023). For EURUS-70B, we mix 63K UltraChat, 30K ShareGPT, and 70K OpenOrca. For LLAMA-3-EURUS-8B, we mix 50k UltraChat, 18k OpenChat (Wang et al., 2023a), TheoremQA decontaminated OpenHermes<sup>4</sup>, and the entire ChatQA (Liu et al., 2024) and CodeAct (Wang et al., 2024).

<sup>2</sup><https://github.com/bigcode-project/bigcode-dataset/tree/main/decontamination>

<sup>3</sup>[https://huggingface.co/datasets/openchat/openchat\\_sharegpt4\\_dataset](https://huggingface.co/datasets/openchat/openchat_sharegpt4_dataset)

<sup>4</sup><https://huggingface.co/datasets/teknium/OpenHermes-2.5>

**Preference Learning.** For hyperparameters, all  $\beta$  is set to 0.1, and  $\lambda_+/\lambda_-$  in KTO is set to 1.33 as recommended. We finetune models for 1 epoch with a 5e-7 learning rate and 0.1 warmup ratio using a cosine scheduler.

**Reward Modeling.** We train RM for 1 epoch with lr=1e-5 learning rate. We also use a cosine scheduler with a warmup ratio of 0.1.

Regarding pair augmentation, we scale up the pairs by matching every correct action for each instruction with one incorrect action of other turns. This leads to NxN pairs of single-turn actions for a trajectory of depth N. We remove the action pairs consisting of nodes at the same turn, as they are already part of the multi-turn trajectory pairs we included. Next, to avoid overfitting on the training set, we only select instructions with  $N \times N \leq 10$ , and for these instructions, we randomly sample at most 9 pairs with each action occurring no more than 3 times. This leads to an augmentation of 240k single-turn action pairs.

## C ADDITIONAL EVALUATION RESULTS OF EURUS

**Detailed Setup in §4.** For math, we test both textual reasoning and program-enhanced settings and report the best performance of the two. All evaluations are conducted in 0-shot CoT with two exceptions: BBH uses 3 shots and IFEval does not use CoT. For MINT, we select MATH, TheoremQA, and MMLU-math from “reasoning” as a new “math” split. We also evaluate 5-shot MMLU (Hendrycks et al., 2021a) for STEM knowledge and MT-Bench (Zheng et al., 2023) for conversation abilities to study whether EURUS needs to trade off other capabilities for reasoning.

**Results.** Results are shown in Table 10.

On MMLU, EURUS outperforms baselines dedicated to coding and math, and achieves comparable or higher results than Mistral-Instruct-v0.2 and Mixtral-8x22B-Instruct-v0.1, the official aligned versions of our base model built by their authors. Compared to general-purpose baseline models, EURUS-7B and LLAMA-3-EURUS-8B achieve comparable performance with the top-performance OpenChat and Starling-LM, and EURUX-8x22B also matches the top level of performance among other general-purpose models.

On MT-Bench, we report baseline numbers from the official leaderboard<sup>5</sup> if available. EURUS matches the performance of mainstream open-source general-purpose models, and EURUX-8x22B further surpasses the score of GPT-3.5 Turbo.

## D DETAILED RESULTS ON REWARD MODELING

### D.1 ADDITIONAL RESULTS ON RERANKING

We present the full results on reranking in Table 11, where the conclusions are consistent with those drawn from §D: (1) Our reward models always achieve the highest accuracy on all test sets across different N, except when N=2 on HumanEval. (2) Both  $\mathcal{L}_{BT}$  and  $\mathcal{L}_{DR}$  consistently help improve reranking performance on three test sets except for HumanEval, where removing either of the objectives can prevent the accuracy from dropping when increasing N from 8 to 16. (3) Modeling

Table 10: MMLU and MT-Bench.

Model	MMLU	MT-Bench
~7B		
Mistral-7B-Instruct-v0.2	58.9	7.60
Zephyr-7B- $\beta$	59.7	7.34
OpenChat-3.5-1210	63.4	7.81
Starling-LM-7B- $\alpha$	64.0	<b>8.09</b>
Magicode-S-DS-6.7B	37.1	4.21
OpenCI-DS-6.7B	37.2	4.06
MAmmoTH-7B-Mistral	56.2	4.25
WizardMath-7B-v1.1	60.3	5.62
OpenMath-Mistral-7B	58.3	2.69
EURUS-7B-SFT	61.8	7.15
+ DPO	62.4	7.38
+ KTO	62.2	7.38
+ NCA	62.2	7.38
LLAMA-3-EURUS-8B-SFT	64.6	6.82
+ DPO	64.8	7.44
+ KTO	<b>64.9</b>	7.27
+ NCA	64.7	7.14
~40B		
Mixtral-8x7B-Instruct	<b>70.3</b>	<b>8.30</b>
DeepSeek-Coder-33B-Ins	40.2	3.83
~70B		
CodeLLaMA-70B-Instruct	55.1	5.12
DeepSeek-LM-67B-Chat	72.3	8.08
QWen1.5-72B-Chat	72.9	8.61
OpenCI-CL-70B	52.4	5.67
OpenMath-CL-70B	60.2	2.29
WizardLM-2-8x22B	77.0	<b>9.08</b>
Mixtral-8x22B-Instruct-v0.1	<b>77.6</b>	8.66
EURUX-8x22B-SFT	75.9	8.20
+ KTO	75.9	8.58
+ NCA	75.6	8.46
Proprietary Models		
GPT-3.5 Turbo	70.0	7.94
GPT-4	<b>86.4</b>	<b>8.96</b>

<sup>5</sup><https://huggingface.co/spaces/lmsys/chatbot-arena-leaderboard>

Table 11: Detailed results of reranking Mistral-Instruct-v0.2’s responses on coding and math.

Datasets	HumanEval				MBPP				GSM8K				MATH			
	2	4	8	16	2	4	8	16	2	4	8	16	2	4	8	16
N																
Random	41.5	39.0	40.2	39.6	33.1	33.6	34.3	30.1	45.0	43.1	44.5	40.2	11.5	11.3	10.0	8.5
Top Logits	43.3	43.3	43.3	43.3	35.3	35.3	35.3	35.3	45.7	45.7	45.7	45.7	12.1	12.1	12.1	12.1
Self-Consistency	43.3	42.7	42.1	40.9	35.3	36.3	36.6	37.1	45.7	49.5	52.2	52.8	12.1	13.8	15.8	16.8
Starling-RM-34B	<b>47.6</b>	<b>47.0</b>	<b>49.4</b>	45.7	37.8	38.8	39.6	40.4	49.1	52.8	56.0	56.5	6.5	7.2	7.7	7.7
EURUS-RM-7B	44.5	45.7	47.6	47.0	39.3	<b>42.6</b>	<b>43.4</b>	<b>43.9</b>	<b>49.8</b>	53.7	56.3	57.3	14.3	16.2	17.1	17.3
w/o $\mathcal{L}_{DR}$	45.7	44.5	46.3	50.0	39.3	42.4	42.4	42.1	49.4	53.2	55.4	56.3	14.2	16.1	17.0	16.9
w/o $\mathcal{L}_{BT}$	45.1	44.5	47.0	48.2	38.6	40.6	39.6	40.1	49.1	52.5	55.2	57.8	14.3	16.3	17.2	17.1
w/o US	45.7	<b>47.0</b>	<b>49.4</b>	<b>50.6</b>	<b>39.3</b>	41.1	41.4	42.9	49.4	<b>53.8</b>	<b>57.4</b>	<b>58.7</b>	<b>14.5</b>	<b>16.6</b>	17.2	<b>17.5</b>
w/o UF + US	43.9	43.3	47.0	46.3	36.3	38.1	36.6	35.3	49.4	52.3	54.6	57.2	14.3	16.5	<b>17.4</b>	17.4
Pass@N	62.8	73.8	88.4	92.7	42.4	48.1	52.6	58.6	54.9	64.1	73.2	80.4	16.9	22.7	28.9	35.5

Table 12: Ablation Study.

Model	Coding			Math					Reasoning	Ins-Following	Avg.
	HumanEval	MBPP	LeetCode	GSM8K	MATH	TheoremQA	SVAMP	ASDiv	BBH	IFEval	
EURUS-7B-SFT	<b>55.5</b>	<b>59.1</b>	<b>20.0</b>	<b>73.7</b>	<b>32.6</b>	20.0	<b>82.2</b>	<b>84.1</b>	64.6	<b>44.0</b>	<b>53.6</b>
Ground-Truth	46.3	46.4	8.9	62.2	15.0	9.6	75.1	68.8	64.4	42.9	44.0
Existing Data Only	38.4	44.1	11.1	45.3	10.8	9.3	52.7	49.4	65.3	43.6	37.0
ULTRAINTERACT Only	46.3	50.1	15.6	67.6	30.9	<b>20.1</b>	80.4	82.0	<b>67.0</b>	17.4	47.7

safety hurts reranking performance in reasoning. When removing UltraSafety from the training data, the RM achieves higher accuracies than EURUS-RM-7B except on MBPP.

## E DETAILED ABLATION RESULTS

We present the full results of §7 in Table 12, with detailed metrics on all coding and math datasets.

## F STATISTICS OF ULTRAINTERACT-V2

We use EURUS-7B-KTO, EURUX-8x22B-NCA, Llama-3-8/70B-Instruct, Llama-3.1-8/70B-Instruct, DeepSeek-Chat-V2 Zhu et al. (2024), DeepSeek-Coder-V2 (Guo et al., 2024a), DeepSeek-Math-RL Shao et al. (2024a), etc, to construct the ULTRAINTERACT-v2. Particularly, the proportion of EURUS responses is 26.4% in the SFT split and 28.2% in the preference split, while the Llama responses take the share of 59.0% and 56.7% in the SFT split and preference split respectively. We present the detailed statistics of ULTRAINTERACT-v2 in Table 13 and Table 14.

## G DATA EXAMPLE

We present an example for SFT data in Table 15 and an example for preference learning in Table 16.

## H ADDITIONAL ABLATION EXPERIMENT RESULTS

We present the ablation study on preference data mixture in this section. We perform KTO on Llama-3-Eurus-8B-SFT with different combinations of UltraFeedback and UltraInteract and the results are presented in Table 17. From the results, we see that training only on UltraInteract leads to higher overall reasoning performances, which are mainly credited to the multi-turn interaction ability, demonstrating the superiority of the tree structure of our data. However, we also observe a lower MT-Bench score compared to training solely on UltraFeedback. Nevertheless, this can be mitigated without hurting reasoning performances by mixing these two datasets together, which indicates that our data is compatible with other datasets, consistent with our conclusions on reward modeling.

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

Table 13: Some statistics of ULTRAINTERACT-v2.

Task	Type		# Instructions	# Turns per Traj.					# Tokens per Traj.	Avg. # Traj per Ins.	Total # Pairs	# Correct Answers
	w/ Interaction?	w/ Tool?		T1	T2	T3	T4	T5				
Math	✓	✓	20,935	8,018	6,149	2,413	1,091	3,264	1,374.2	1.0	61,738	80,982
	✗	✓	1,279	6,723	-	-	-	-	423.0	5.3	2,055	6,723
	✓	✗	21,474	10,117	4,707	1,263	589	4,798	1,155.7	1.0	72,944	143,587
	✗	✗	634	4,876	-	-	-	-	396.0	7.7	2,172	4,876
Coding	✓	-	9,498	4,323	2,057	942	459	1,717	1,429.4	1.0	71,900	155,429
	✗	-	2,877	9,824	-	-	-	-	591.0	3.4	5,990	9,824
<b>Total</b>	-	-	56,697	43,881	12,913	4,618	2,139	9,779	894.9	3.2	216,799	410,421

Table 14: Stats breakdown for ULTRAINTERACT-v2

Task	Dataset	w/ Tool?	# Prompts	# Pairs	# Correct Answers.	Avg. Length	Human Annotation	
							Has Answer?	Has Rationale?
Math	GSM8K	✓	4,438	14,616	17,797	642.4	✓	✓
		✗	6,997	10,708	11,893	367.6	✓	✓
	MATH	✓	6,197	20,056	27,217	681.8	✓	✓
		✗	6,799	46,720	115,924	531.7	✓	✓
	MathQA	✓	6,046	19,315	22,062	653.8	✓	✓
		✗	5,898	14,525	16,237	570.3	✓	✓
	NumGLUE	✓	2,436	4,241	9,113	570.4	✓	✗
		✗	2,414	3,163	4,409	409.0	✓	✗
TabMWP	✓	3,097	5,565	11,516	507.1	✓	✗	
Coding	CodeContest	-	5,788	39,713	76,272	1,291.7	✓	✓
	TACO	-	6,459	37,317	86,844	1,062.9	✓	✓
	Codeforces	-	128	860	2,137	1,480.9	✓	✓

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

Table 15: An example from the SFT data.

**Task: Coding**

**Instruction:** Write Python code to solve the task.

Mike has a string  $s$  consisting of only lowercase English letters. He wants to change exactly one character from the string so that the resulting one is a palindrome.

A palindrome is a string that reads the same backward as forward, for example strings  $z$ ,  $aaa$ ,  $aba$ ,  $abccba$  are palindromes, but strings  $codeforces$ ,  $reality$ ,  $ab$  are not.

Input

The first and single line contains string  $s$ .

Output

Print YES (without quotes) if Mike can change exactly one character so that the resulting string is palindrome or NO (without quotes) otherwise.

**Response:**

Step 1: This problem asks us to determine if it is possible to change exactly one character in the string to make it a palindrome.

Step 2: We can start by iterating through the characters in the string and counting the number of differences between corresponding characters from the beginning and end of the string.

Step 3: If the length of the string is odd and there is only one difference, or if the length is even and there is exactly one difference, then it is possible to change exactly one character to make the string a palindrome.

Step 4: We need to consider edge cases such as a single character string or an empty string. `python`

# Step 5: Take the input string

`s = input()`

# Step 6: Calculate the length of the string

`n = len(s)`

# Step 7: Initialize a counter to keep track of the number of differences

`cnt = 0`

# Step 8: Iterate through the first half of the string and count the number of differences between corresponding characters from the beginning and end of the string.

`for i in range(n//2):`

`if s[i] != s[n-i-1]:`

`cnt += 1`

# Step 9: If the length of the string is odd and there are no differences, or if the length is even and there is exactly one difference, print "YES" (without quotes)

`if (n % 2 == 1 and cnt == 0) or (cnt == 1):`

`print('YES')`

# Step 10: Otherwise, print "NO" (without quotes)

`else:`

`print('NO')`

Table 16: An example from the Preference data.

1134	1135	1136	1137	1138	1139	1140	1141	1142	1143	1144	1145	1146	1147	1148	1149	1150	1151	1152	1153	1154	1155	1156	1157	1158	1159	1160	1161	1162	1163	1164	1165	1166	1167	1168	1169	1170	1171	1172	1173	1174	1175	1176	1177	1178	1179	1180	1181	1182	1183	1184	1185	1186	1187
<b>Task:Math</b>																																																					
<b>Trajectory:</b>																																																					
[ {																																																					
from: user,																																																					
value: Tool available:[1] Python interpreter. When you send a message containing Python code to python, it will be executed in a stateful Jupyter notebook environment. Solve the following math problem step-by-step. Simplify your answer as much as possible. Simplify $(-3 - 2i) - (1 + 4i)$ .																																																					
},																																																					
{																																																					
from: assistant,																																																					
value: Step 1: To simplify the expression $(-3 - 2i) - (1 + 4i)$ , we need to distribute the negative sign to both terms inside the parentheses. Step 2: We can then combine like terms by adding the real parts and the imaginary parts separately. # Step 3: Distribute the negative sign to both terms inside the parentheses. $(-3) - (-1) + (-2i) - (-4i)$ # Step 4: Combine like terms by adding the real parts and the imaginary parts separately. $-3 + 1 + (-2i + 4i)$																																																					
},																																																					
{																																																					
from: user,																																																					
value: Out: Cell In[1], line 2 $(-3) - (-1) + (-2i) - (-4i)$ SyntaxError: invalid syntax. Your answer is wrong. Your thought process for simplifying the given complex numbers is correct. However, you made a syntax error in the code. In Python, the multiplication operation must be explicitly written out with an asterisk (*), even for complex numbers. Remember, the imaginary unit in Python is denoted as j and not i. Try correcting the syntax and calculate again.																																																					
},																																																					
{																																																					
from: assistant,																																																					
value: Step 1: To simplify the expression $(-3 - 2i) - (1 + 4i)$ , we need to distribute the negative sign to both terms inside the parentheses. Step 2: We can then combine like terms by adding the real parts and the imaginary parts separately. # Step 3: Distribute the negative sign to both terms inside the parentheses. $(-3) - (-1) + (-2j) - (-4j)$ # Step 4: Combine like terms by adding the real parts and the imaginary parts separately. $-3 + 1 + (-2j + 4j)$																																																					
},																																																					
{																																																					
from: user,																																																					
value: Out: $(-2+2j)$ . Your answer is wrong. Your thought process for simplifying the given complex numbers is on the right track. However, you made an arithmetic error when adding the real parts and the imaginary parts. You should have combined the like terms correctly to get the right answer $-4$ and $-6j$ respectively. Please try again with the correct arithmetic calculations.																																																					
} ]																																																					
<b>Chosen:</b>																																																					
Step 1: Split the expression into real and imaginary terms.																																																					
Step 2: Subtract the real terms from each other.																																																					
Step 3: Subtract the imaginary terms from each other.																																																					
Step 4: Combine the results of the subtraction to get the simplified form of the expression.																																																					
# Step 5: Split the expression into real and imaginary parts and then subtract them.																																																					
real_part = -3 - 1 parts																																																					
imaginary_part = -2 - 4																																																					
# Step 6: Combine the real and imaginary parts to get the result.																																																					
result = f{real_part}{imaginary_part:+d}i																																																					
print(result)																																																					
<b>Rejected:</b>																																																					
Step 1: To simplify the expression $(-3 - 2i) - (1 + 4i)$ , we need to distribute the negative sign to both terms inside the parentheses.																																																					
Step 2: We can then combine like terms by adding the real parts and the imaginary parts separately.																																																					
# Step 3: Distribute the negative sign to both terms inside the parentheses.																																																					
$(-3) - 1 - (-2j) - 4j$																																																					
# Step 4: Combine like terms by adding the real parts and the imaginary parts separately.																																																					
$-3 - 1 - 2j - 4j$																																																					

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

Table 17: Additional experiment results on preference learning with different compositions of UltraFeedback and UltraInteract.

<b>Model</b>	<b>Coding</b>	<b>Math</b>	<b>BBH</b>	<b>IFEval</b>	<b>Multi-Turn</b>	<b>Avg.</b>	<b>MT-Bench</b>
<b>Llama-3-Eurus-8B-SFT</b>	42.11	53.97	72.40	47.10	21.45	46.60	6.8
<b>+KTO (UF + UI)</b>	41.84	56.14	71.70	50.60	31.92	49.40	7.3
<b>+KTO (Only UF)</b>	44.60	56.26	72.00	50.10	24.03	48.80	7.5
<b>+KTO (Only UI)</b>	40.70	55.86	71.70	50.60	34.49	49.40	7.2