

Adaptive Preference Optimization with Uncertainty-aware Utility Anchor

Xiaobo Wang^{1,3}, Zixia Jia³, Jiaqi Li³, Qi Liu^{1,2*}, Zilong Zheng^{3*}

¹State Key Laboratory of Cognitive Intelligence, University of Science and Technology of China

²Institute of Artificial Intelligence, Hefei Comprehensive National Science Center

³State Key Laboratory of General Artificial Intelligence, BIGAI

xiaobow.0519@gmail.com, qiliuql@ustc.edu.cn,

{lijiaqi, jiazixia, zlzheng}@bigai.ai

Abstract

Offline preference optimization methods are efficient for large language models (LLMs) alignment. Direct Preference optimization (DPO)-like learning, one of the most popular approaches, stands out for its efficiency in reward modeling. However, these methods typically follow the convention to use Bradley-Terry (BT) reward modeling that faces several critical assumptions, including the requirement for pairwise training data, model distribution shifting, human rationality assumption, etc. To address these limitations, we propose a general framework for offline preference optimization methods, Adaptive Preference Optimization with Utility Anchor (UAPO), which introduces an anchoring function to estimate the uncertainties brought from preference data annotation. Our method enables training even in scenarios where the data is unpaired, significantly enhancing data utilization efficiency. Moreover, the anchor design makes UAPO more robust in the training process. Experimental results demonstrate that UAPO achieves competitive outcomes without the strict dependency on data pairing, paving the way for more flexible and effective preference optimization methods.

1 Introduction

Aligning large language models (LLMs) with human values and preferences is a crucial step in LLM training. This alignment ensures that LLMs can better follow human instructions, becoming more helpful, harmless, and honest. However, ensuring that these models align with human preferences and ethical standards remains a significant hurdle. Previous work has made strides in this direction by employing techniques. Reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) is a widely adopted method in the alignment domain, yielding significant improvements

in LLM performance. However, the high computational cost and time consumption associated with RLHF present substantial challenges, limiting its broader application. To address these challenges, offline preference optimization methods, such as DPO (Rafailov et al., 2023), have been developed as a more manageable offline training process. This shift makes the training procedure more controllable and resource-efficient.

Despite the inspiring progress and superior benchmark performance made, current methodologies predominantly adhere to the Bradley-Terry (BT; Bradley and Terry (1952)) model or its generalized form, the Plackett-Luce model (Plackett, 1975), for reward modeling, a rank-based method widely applied to predict rational preference distribution given pairwise data. In this work, we challenge this convention in Section 2.2 by thoroughly analysis of its discrepancies with real-world scenarios. From the **data** perspective, the BT model’s reliance on pairwise comparisons imposes a significant constraint on the collection of preference data, such as identifying winning and losing responses from a superior model or a pre-trained reward function. This process is both labor-intensive and data-inefficient. From the **optimization** perspective, over-optimization in preference learning, as highlighted by Goodhart’s law (Gao et al., 2023), can lead to distribution shifts and reward hacking when applied to out-of-distribution samples. This may also lead to a disparity between training and testing reward functions. From the **cognitive** perspective, the BT model presupposes that human annotators are fully rational and maximize utility values, which has been challenged by behavioral economics (Tversky and Kahneman, 1992) that humans are typically risk-averse. These limitations are also observed in synthetic data (Cui et al., 2023) generated by modern language models.

In this work, we introduce Adaptive Preference Optimization with Uncertainty-aware Utility An-

*Corresponding authors

chor (UAPO), a cognitively inspired preference optimization framework designed to address the above limitations across data, learning, and cognitive dimensions. Specifically, drawing inspiration from the anchoring effect identified in behavioral economics by Simon (1955), UAPO introduces the concept of a *utility anchor* to better model human preferences in a more flexible and robust manner to solve the uncertainty in the data labeling process. This integration allows for more efficient processing of unpaired data and yields a more precise estimation of human preferences. By decoupling the learning objective into unpaired forms, the utility anchor enhances robustness to various hyperparameter settings, eliminating the necessity for meticulous customization or tuning across different models and datasets. We show that UAPO is flexible and can easily be integrated to enhance previous offline reward learning methods, such as DPO (Rafailov et al., 2023) and SimPO (Meng et al., 2024). Additionally, we provide theoretical analysis from the perspective of uncertainty penalty, a strategy commonly used in pessimistic RL that use uncertainty rewards (Jin et al., 2021) to address the overoptimization problem in offline reward learning.

In Section 3, we run extensive experiments to demonstrate the effectiveness of our UAPO and its variant SimUAPO. We highlight the main advantages of UAPO as follows.

- **A general framework for offline preference optimization:** UAPO can transform most offline preference optimization methods, which typically require pairwise data, into methods capable of effectively utilizing unpaired data.
- **The utility anchor helps mitigate uncertainty in data labeling:** By naturally leveraging unpaired data, the utility anchor is inherently aware of uncertainty. In UAPO, it further guides the model to effectively handle ambiguous or inconsistent preference signals.
- **Rigorous theoretical analysis of utility anchor:** we prove that the utility anchor is more consistent with concurrent offline preference modeling, making the alignment process more robust.

Our results highlight the potential of UAPO as a versatile and robust solution for preference optimization, paving the way for its application in real-world scenarios where data and computational resources are limited.

2 Uncertainty-aware Utility Anchor

In this section, we start by briefly discussing prior offline preference learning methods (§2.1). Then we thoroughly analyze critical assumptions of reward modeling in previous methods that do not hold in real-world settings (§2.2). To address the issues, we introduce a cognitive compelling utility anchor and devise our formulation of UAPO in §2.3. We further provide theoretical justification in §2.4 from the perspective of pessimistic RL.

2.1 Background: Offline Preference Learning

Consider an instructed preference dataset $\mathcal{D} = \{x, y_w, y_l\}$, where x represents the prompt, and (y_w, y_l) are two responses generated by a reference model π_{ref} . y_w signifies the preferred (winning) response, and y_l denotes the dispreferred (losing) response, both of which can be annotated by human annotators or a language model. The function $r(x, y)$ is a latent reward function employed to model the preference comparison between y_w and y_l using the Bradley-Terry (BT) model (Bradley and Terry, 1952), or its generalized form Plackett-Luce (Plackett, 1975), given by:

$$p(y_w \succ y_l | x) = \frac{\exp(r(x, y_w))}{\exp(r(x, y_w)) + \exp(r(x, y_l))} = \sigma(r(x, y_w) - r(x, y_l)), \quad (1)$$

where $\sigma(\cdot)$ denotes the sigmoid function.

Direct Preference Optimization (DPO) (Rafailov et al., 2023) is a leading offline preference optimization method that reparameterizes the reward function $r(x, y)$ using a closed-form expression with the optimal policy:

$$r_{\text{DPO}}(x, y) = \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} + \beta \log Z(x), \quad (2)$$

where π_{θ} is the policy model, π_{ref} is the reference policy (typically the SFT model), and $Z(x)$ is the partition function. Integrating $r_{\text{DPO}}(x, y)$ into BT model yields the loss for DPO:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]. \quad (3)$$

Simple Preference Optimization (SimPO) (Meng et al., 2024), as a variant of DPO, demonstrates superior efficacy across diverse tasks by utilizing a reference-free reward with response length normalization:

$$r_{\text{SimPO}}(x, y) = \frac{\beta}{|y|} \log \pi_{\theta}(y|x), \quad (4)$$

Method	Objective
DPO	$-\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right)$
IPO	$\left(\log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} - \frac{1}{2\tau} \right)^2$
CPO	$-\log \sigma \left(\beta \log \pi_{\theta}(y_w x) - \beta \log \pi_{\theta}(y_l x) \right) - \lambda \log \pi_{\theta}(y_w x)$
KTO	$-\lambda_w \sigma \left(\beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - z_0 \right) + \lambda_l \sigma \left(z_0 - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right)$, where $z_0 = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\beta \text{KL}(\pi_{\theta}(y x) \parallel \pi_{\text{ref}}(y x))]$
ORPO	$-\log p_{\theta}(y_w x) - \lambda \log \sigma \left(\frac{p_{\theta}(y_w x)}{1-p_{\theta}(y_w x)} \right) - \log \left(\frac{p_{\theta}(y_l x)}{1-p_{\theta}(y_l x)} \right)$, where $p_{\theta}(y x) = \exp \left(\frac{1}{ y } \log \pi_{\theta}(y x) \right)$
R-DPO	$-\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} + (\alpha y_w - \alpha y_l) \right)$
SimPO	$-\log \sigma \left(\frac{\beta}{ y_w } \log \pi_{\theta}(y_w x) - \frac{\beta}{ y_l } \log \pi_{\theta}(y_l x) - \gamma \right)$
UAPO	$-\log \sigma \left(r(x, y_w) - r(x, y_{\perp}) \right) - \log \sigma \left(r(x, y_{\perp}) - r(x, y_l) \right)$

Table 1: Comparison of different offline preference optimization methods and their objectives.

where $|y|$ is the length of response y . A target reward margin γ is introduced to distinguish between the preference rewards ensuring the reward difference between winning and losing responses exceeds this margin. The loss of SimPO is as below:

$$\mathcal{L}_{\text{SimPO}} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(r_{\text{SimPO}}(x, y_w) - r_{\text{SimPO}}(x, y_l) - \gamma \right) \right]. \quad (5)$$

2.2 Limitations of offline preference optimization methods

In this study, we delineate three critical issues from the perspectives of data, optimization, and cognition, that are typically treated as conventional in previous formulations of preference modeling (Table 1).

Pairwise Training Data From the **data** perspective, most preference optimization methods rely on pairwise comparison data to align a policy model π_{θ} , requiring multiple responses to be compared to determine alignment with human preferences. However, human preferences are often expressed without comparisons, such as *stating a favorite fruit without indicating dislikes*. Such rationale also follows prior literature in prospect theory (Tversky and Kahneman, 1992), later formalized as KTO (Ethayarajh et al., 2024), claiming that the reward function shall be a human-aware objective, in which they incorporate a pre-defined middle state z_0 to estimate the expected KL distance between $\pi_{\theta}(x_{\text{data}})$ and $\pi_{\text{ref}}(x_{\text{data}})$, which is nevertheless computationally intractable. Therefore, this estimation can only be approximated through sampling, which often leads to suboptimal performance.

Distribution Shifting From the **optimization** perspective, it has been well studied the issue of over-optimization (Gao et al., 2023) in preference learning. In accordance with Goodhard’s law, optimizing reward functions with specific in-distribution training data can inevitably lead to a distribution shift of outputs given an out-of-distribution (OOD) sample, *i.e.*, reward hacking. In this context, the OOD sample corresponds to real-world human preferences that deviate from those represented in the training distribution. Moreover, a discrepancy is created between training and testing reward functions: for training, a relative reward *w.r.t.* losing sample y_l is learned; and for testing, the policy π_{θ} has to generate a response solely based on input x without relying on y_l and reference model π_{ref} . This mismatch can cause the reward model to produce unreliable signals, thereby impairing the effectiveness of the policy during inference.

“Irrational” Preference Pairs From the **cognitive** perspective, a critical assumption of the Bradley-Terry (BT) model (Bradley and Terry (1952); Eq. (1)) is that annotators behave as fully rational agents who make pairwise comparisons by maximizing utility values with respect to an absolute, universal utility function (Fisch et al., 2024; Azar et al., 2024). However, this assumption often fails to hold in practice, as irrational behaviors are commonly observed, both in the preference data annotation process, where annotator judgments may be inconsistent or biased, and in the optimization stage, where model updates may not strictly follow utility-maximizing principles. Moreover, behavioral economists have identified that humans are *risk averse* (Tversky and Kahneman, 1992): decisions are made in *relative* terms rather than absolutes. This phenomenon is also observed in modern language models (Jia et al., 2024), which are frequently employed as automated preference data generators (Cui et al., 2023) or evaluation judges (Li et al., 2024b). Such “irrationality” results in the modeling of moderate rewards rather than true utility functions, as presented in our analysis of reward comparisons (Figure 1(a)).

2.3 Reward Modeling with Utility Anchors

One critical challenge of the above issues is that the true distribution $p_{\text{human}}(y_w \succ y_l|x)$ is not known, while the observed data pairs $(x, y_w, y_l) \sim \mathcal{D}$ are subject to human aware uncertainties. In this section, we borrow the idea of “reference anchors” (Tversky and Kahneman, 1992) and propose

a more robust and flexible preference optimization framework in terms of data efficiency, reward shifting, and cognitive compelling.

Anchoring Effect The anchoring effect, originally raised by behavior economist [Simon \(1955\)](#), describes the common human tendency to rely heavily on a reference datapoint (the "anchor"), mostly the first information, when making decisions. Inspired by the cognitive bias, we curate a **utility anchor** that is better in line with concurrent offline preference modeling, defined as:

$$r(x, y_{\perp}) = \beta \log \frac{\pi_{\theta}(x, y_{\perp})}{\pi_{\text{ref}}(x, y_{\perp})} + \beta \log Z(x) + \gamma, \quad (6)$$

where y_{\perp} is a dummy token that can be learned through preference optimization, γ is a constant reward anchor (more details can be seen in [Appendix C](#)). By introducing the utility anchor, the model gains the ability to handle uncertainty in preference data labeling ([§2.2](#)), allowing it to deal with ambiguous or inconsistent signals more effectively.

Since the sigmoid function is monotonically increasing, we have

$$\begin{aligned} \log \sigma(r(x, y_w) - r(x, y_{\perp})) + \log \sigma(r(x, y_{\perp}) - r(x, y_l)) \\ \leq \log \sigma(r(x, y_w) - r(x, y_l)). \end{aligned} \quad (7)$$

Thus, given the definition of utility anchor, we can rewrite the log-likelihood form of [Eq. \(1\)](#) into

$$\begin{aligned} \log p(y_w \succ y_l | x) = \log \sigma(r(x, y_w) - r(x, y_{\perp})) \\ + \log \sigma(r(x, y_{\perp}) - r(x, y_l)). \end{aligned} \quad (8)$$

The UAPO Objective Finally, we obtain the UAPO objective:

$$\begin{aligned} \mathcal{L}_{\text{UAPO}}(\pi_{\theta}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma(r(x, y_w) - r(x, y_{\perp})) \right. \\ \left. + \log \sigma(r(x, y_{\perp}) - r(x, y_l)) \right], \end{aligned} \quad (9)$$

where $r(x, y)$ denotes the reward function in any offline preference optimization method, and $r(x, y_{\perp})$ represents the utility anchor corresponding to the adopted method. For example, when applying SimPO within this framework, the reference-free reward $r_{\text{SimUAPO}}(x, y)$ is given in [Eq. \(4\)](#), and the utility anchor can be reformulated as:

$$r_{\text{SimUAPO}}(x, y_{\perp}) = \beta \log \pi_{\theta}(x, y_{\perp}) + \gamma. \quad (10)$$

Learning from unpaired data Most preference datasets exhibit a strong win/lose imbalance, where dispreferred responses significantly outnumber preferred ones. For example, in UltraFeedback ([Cui](#)

[et al., 2023](#)), each prompt is used to query multiple LLMs to generate four candidate responses, which are then annotated by GPT-4. Following previous works, the responses are categorized into one winning response and three losing responses, resulting in data imbalance. To better solve the data imbalance problem and harness the capability of UAPO to effectively process unpaired data, we construct multiple datasets ([§3.1](#)). And considering the separate form of [Eq. \(8\)](#), we could expand the objective function to accommodate multiple datasets. Given a set of n responses $\{y_1, \dots, y_n\}$ generated by different language models, we aim to induce an utility anchor hidden within human (or LLM) judges, *s.t.* $(y_{w1} \succ y_{l1}), (y_{w2} \succ y_{l2}), \dots$. Then the above objective turns into (refer to [Appendix A](#) for proof):

$$\begin{aligned} \mathcal{L}_w &= -\sum_{y_w} \log \left(\frac{\exp(r(x, y_w))}{\exp(r(x, y_{\perp})) + \sum_{y'_w} \exp(r(x, y'_w))} \right) \\ \mathcal{L}_l &= -\log \left(\frac{\exp(r(x, y_{\perp}))}{\exp(r(x, y_{\perp})) + \sum_{y_l} \exp(r(x, y_l))} \right) \\ \mathcal{L}_{\text{UAPO}_{\text{multi}}} &= \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} (\mathcal{L}_w + \mathcal{L}_l). \end{aligned} \quad (11)$$

This adjustment allows UAPO to effectively exploit the availability of preferred and dispreferred responses without the need for explicitly paired counterparts.

Among the methods in [Table 1](#), KTO also supports learning from unpaired data but is based on prospect theory, while UAPO is motivated by uncertainty penalties in reinforcement learning optimization. KTO relies on a fixed middle state z_0 to estimate the expected KL divergence between $\pi_{\theta}(x_{\text{data}})$ and $\pi_{\text{ref}}(x_{\text{data}})$, which is approximated via prompt sampling and often underperforms. In contrast, UAPO learns a utility anchor through a trainable linear representation, offering better efficiency and performance.

2.4 Theoretical Analysis: Preference Optimization with Uncertainty Penalty

From the optimization objective, we provide further interpretation of the UAPO from the view of pessimistic RL, and justify how it addresses the distribution shifting given the OOD samples ([§2.2](#)).

The utility anchor defined in [Eq. \(6\)](#) can be rewritten into an output-independent score:

$$r_{\phi}(x) = \beta \log \frac{u_{\phi}(\pi_{\theta}(x))}{u_{\phi}(\pi_{\text{ref}}(x))} + \beta \log Z(x) + \gamma, \quad (12)$$

where u is a linear function parameterized by ϕ to model the utility anchor y_{\perp} . Applying this to the

UAPO objective in Eq. (9) and set $\gamma = 0$ for simplicity and fair theoretical comparison with DPO:

$$\begin{aligned} \max_{\pi_\theta, u_\phi} \mathbb{E} & \left[\sigma \left(\beta \log \frac{\pi_\theta(x, y_w)}{\pi_{\text{ref}}(x, y_w)} - \beta \log \frac{u_\phi(\pi_\theta(x))}{u_\phi(\pi_{\text{ref}}(x))} \right) \right. \\ & \left. + \sigma \left(\beta \log \frac{u_\phi(\pi_\theta(x))}{u_\phi(\pi_{\text{ref}}(x))} - \beta \log \frac{\pi_\theta(x, y_l)}{\pi_{\text{ref}}(x, y_l)} \right) \right]. \end{aligned} \quad (13)$$

Given a learned u_ϕ , the derivative of the preferred objective \mathcal{L}_w (we skip the symmetric form of \mathcal{L}_l for simplicity) is

$$\begin{aligned} \nabla \mathcal{L}_w = & -\beta \mathbb{E} \left[\sigma \left(\underbrace{-\left(\hat{r}_\theta(x, y_w) - \hat{u}_{\phi, \theta}(x) \right)}_{\text{uncertainty normalized reward}} \right) \right. \\ & \left. \left[\nabla_\theta \log \pi(y_w|x) - \nabla_\theta \log u(\pi(x)) \right] \right], \end{aligned} \quad (14)$$

where $\hat{r}_\theta(x, y) = \beta \log \pi(x, y) / \pi_{\text{ref}}(x, y)$ and $\hat{u}(x) = \beta \log u_\phi(\pi_\theta(x)) / u_\phi(\pi_{\text{ref}}(x))$ can be considered a **uncertainty reward penalty**. Considering the increasing feature of the Sigmoid function, the objective is equivalent to optimizing the lower-bound of an uncertainty-aware policy likelihood (refer to Appendix B for proof):

$$\log \pi_{\theta, \phi}(y|x) = \log \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) e^{\frac{1}{\beta}(r(x, y) - \hat{u}(x))}, \quad (15)$$

where $\hat{u}(x)$ indicates the reference anchor of uncertainties given input x . The representation is also consistent with conventional pessimistic RL theories (Jin et al., 2021) that produce a conservative estimate of reward as a lower confidence bound.

3 Experiments

3.1 Setup

Models and Training Dataset We select the instruct models of Mistral-7B (Jiang et al., 2023), Llama-3-8B (AI@Meta, 2024), and the base model of Gemma-2-9B (Team et al., 2024) for evaluation to acquire better instruction-following capabilities in this paper.

We use three datasets including princeton-nlp/mistral-instruct-ultrafeedback¹, princeton-nlp/llama3-ultrafeedback-armorm², and princeton-nlp/gemma2-ultrafeedback-armorm³ for training Mistral-7B, Llama-3-8B and Gemma-2-9B respectively. It is worth mentioned that we further reconstruct the training data in these three datasets to

¹<https://huggingface.co/datasets/princeton-nlp/mistral-instruct-ultrafeedback>

²<https://huggingface.co/datasets/princeton-nlp/llama3-ultrafeedback-armorm>

³<https://huggingface.co/datasets/princeton-nlp/gemma2-ultrafeedback-armorm>

obtain one winning response accompanied with three losing responses, rather than one pair of winning and losing response for comparison. More details about **multiple datasets** construction can be seen in Appendix D.

Baselines We compare UAPO with various offline preference optimization methods including DPO (Rafailov et al., 2023), IPO (Azar et al., 2024), CPO (Xu et al., 2024), KTO (Ethayarajh et al., 2024), ORPO (Hong et al., 2024), R-DPO (Park et al., 2024) and SimPO (Meng et al., 2024). Each method addresses distinct challenges in aligning policy models with human preferences.

Benchmarks and Evaluation Metrics We evaluate our models on three widely recognized open-ended instruction-following benchmarks: AlpacaEval 2 (Li et al., 2023), Arena-Hard v0.1 (Li et al., 2024a,b), and MT-Bench (Zheng et al., 2023). These benchmarks are considered representative of human preference evaluations relative to the training data. AlpacaEval 2 consists of 805 questions sampled from five datasets, covering various instruction-following tasks. MT-Bench includes 80 questions distributed across eight categories, providing a structured evaluation of model performance on different conversational and reasoning tasks. Arena-Hard v0.1 extends MT-Bench by introducing 500 well-defined technical problem-solving questions, further challenging the models' reasoning and problem-solving abilities.

We consistently follow the usage of the evaluation metrics that are used in the original benchmark for each benchmark. For AlpacaEval 2, we present both the raw win rate (WR) and the length-controlled win rate (LC) (Dubois et al., 2024) which is designed to mitigate the influence of model verbosity on preference judgments. For Arena-Hard, we use WR relative to a baseline model, allowing direct comparison of alignment performance. For MT-Bench, we compute the average score using GPT-4 as the judge model for assessing the instruction-following quality.

We also evaluate our models on RewardBench (Lambert et al., 2025) and RewardBench 2 (Malik et al., 2025) to compare performance in the OOD settings. RewardBench and RewardBench 2 are widely known as comprehensive benchmarks for evaluating reward models, specifically focusing on preference alignment and generalization ability beyond training distributions. RewardBench provides a diverse set of prompts and human preference annotations designed to capture realistic deployment

Method	Mistral-7B-Instruct				Llama-3-8B-Instruct				Gemma-2-9B-Instruct			
	AlpacaEval 2		Arena-Hard	MT-Bench	AlpacaEval 2		Arena-Hard	MT-Bench	AlpacaEval 2		Arena-Hard	MT-Bench
	LC (%)	WR (%)	WR (%)	GPT-4	LC (%)	WR (%)	WR (%)	GPT-4	LC (%)	WR (%)	WR (%)	GPT-4
pairwise datasets												
SFT	19.0	15.4	12.9	7.5	26.0	25.3	22.3	8.1	48.7	36.5	42.1	8.6
IPO	20.3	20.3	16.2	7.8	46.8	42.4	36.6	8.2	62.6	58.4	53.5	8.7
CPO	23.8	28.8	22.6	7.5	34.1	36.4	30.9	<u>8.2</u>	56.4	53.4	55.2	8.7
KTO	24.5	23.6	17.9	<u>7.7</u>	34.1	32.1	27.3	<u>8.2</u>	61.7	55.5	53.8	8.5
ORPO	24.5	24.9	<u>20.8</u>	<u>7.7</u>	38.1	33.8	28.2	8.3	56.2	46.7	46.2	8.3
R-DPO	27.3	24.5	<u>16.1</u>	7.5	48.0	<u>45.8</u>	<u>35.1</u>	8.2	68.3	66.9	57.9	8.5
DPO*	26.8	24.9	16.3	7.6	48.2	47.5	35.2	8.2	67.8	65.4	58.9	-
DPO	20.5	18.1	13.4	7.6	41.2	37.5	33.9	8.2	68.8	64.5	58.4	8.7
SimPO*	32.1	34.8	21.0	7.6	53.7	47.5	36.5	8.0	72.4	65.9	57.8	-
SimPO	29.0	<u>31.9</u>	19.9	7.2	<u>50.1</u>	45.2	27.7	7.4	<u>72.6</u>	66.1	57.8	<u>8.8</u>
UAPO	23.1	17.7	11.1	<u>7.7</u>	41.2	38.3	32.1	<u>8.2</u>	70.2	67.7	<u>58.8</u>	8.9
SimUAPO	<u>28.6</u>	32.2	17.7	7.6	51.2	47.6	34.4	8.3	73.5	<u>67.0</u>	59.4	8.9

Table 2: Performance comparison of different methods on Mistral-Instruct, Llama-3-Instruct, and Gemma-2-Instruct trained with pairwise datasets. \star denotes results reported by Meng et al. (2024). The performance differences of DPO and SimPO are attributed to the version of GPT used. The best results are highlighted in **bold**, while the second-best results are underlined.

scenarios, while RewardBench 2 extends this setup with a broader and more challenging collection of tasks that place greater emphasis on consistency and robustness across domains. Together, these two benchmarks capture real-world human preferences and enable us to assess how well our models adapt to distribution shifts between training data and practical human evaluations.

3.2 Main Results

UAPO gains an improvement over existing preference optimization methods on pairwise datasets As shown in Table 2, all offline preference optimization methods outperform the SFT baseline, demonstrating the effectiveness of preference-based training. Among them, UAPO and SimUAPO consistently achieve superior results on Llama-3-Instruct and Gemma-2-Instruct compared to their respective counterparts, DPO and SimPO. Specifically, on AlpacaEval 2, UAPO outperforms DPO by 2.6 points on Mistral-Instruct and 1.4 on Gemma-2-Instruct, showing stronger alignment with human preferences. SimUAPO achieves state-of-the-art LC results on Gemma-2-Instruct. On Arena-Hard, it improves over SimPO by 6.7 points on Llama-3-Instruct and 1.6 on Gemma-2-Instruct.

Meanwhile, we observe that MT-Bench scores are relatively consistent across the three models and all preference optimization methods. This may be due to the inherent randomness in MT-Bench scoring, as discussed in Appendix E.4. While SimUAPO performs comparably to SimPO on Mistral-Instruct and shows a slight drop in LC on

Method	AlpacaEval 2		Arena-Hard	MT-Bench
	LC (%)	WR (%)	WR (%)	GPT-4
multiple datasets				
DPO	41.2	37.5	33.9	8.2
SimPO	50.1	45.2	27.7	7.4
UAPO	41.2	38.3	32.1	8.2
SimUAPO	<u>51.2</u>	47.6	<u>34.4</u>	<u>8.3</u>
DPO-multi	43.3	38.0	30.1	7.2
SimPO-multi	45.7	41.6	27.5	6.7
UAPO-multi	44.0	43.0	33.5	7.3
SimUAPO-multi	55.2	<u>47.2</u>	35.2	8.4

Table 3: Performance comparison of different methods on Llama-3-Instruct trained with multiple datasets.

AlpacaEval 2, this could be attributed to training dynamics specific to Mistral-Instruct.

In general, the consistently strong results across different models and benchmarks show the effectiveness of UAPO in preference-based alignment.

Multiple datasets can improve UAPO much more We utilize the multiple datasets mentioned in Section 3.1 to train DPO, SimPO, UAPO and SimUAPO on Llama-3-Instruct. In Table 3, DPO-multi, SimPO-multi, UAPO-multi, and SimUAPO-multi refer to the respective training results using multiple datasets. And we set constant reward anchor $\gamma = 0$ in UAPO-multi training process. Compared to using pairwise datasets, UAPO and SimUAPO show a significant improvement in training on multiple datasets.

Overall, UAPO and SimUAPO demonstrate superior performance compared to DPO and SimPO when utilizing multiple datasets. Specifically, on Llama-3-Instruct, SimUAPO-multi achieves a 55.2 LC on AlpacaEval 2, surpassing other offline preference optimization methods. However, the MT-

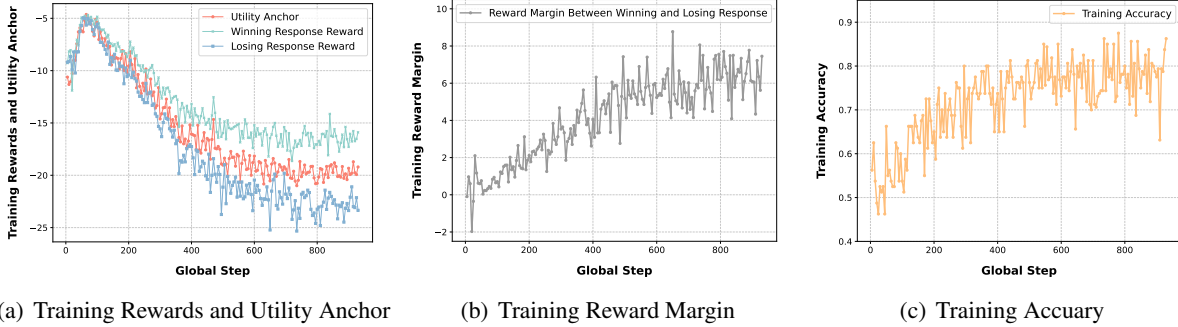


Figure 1: Analysis of utility anchor representation during SimUAPO training on Gemma-2-Instruct: (a) Relationship between the margin of winning and losing response reward; (b) Reward margin during training, and (c) Accuracy for training, i.e., the percentage of response rewards for winning greater than those of losing.

Prompt	Reward of y_w	Utility Anchor	Reward of y_l
Prompt A: <i>Teacher: In this task you will be given a list of integers. You should remove any integer that is not prime ...</i>	-9.875	-10.764	-11.587
Prompt B: <i>explain all machine learning algorithms with examples and code in python</i>	-11.346	-19.879	-21.724

Table 4: The intuitive example cases of Gemma-2-Instruct trained by SimUAPO demonstrate that the utility anchor is closely related to the prompt. The reward of y_w for Prompt B is lower than the utility anchor of Prompt A. Therefore, if a constant standard is used, y_w from Prompt B will be classified as part of the losing responses.

Bench scores remain inconsistent, as their trends do not align with those observed on other benchmarks, further supporting our earlier observation regarding the inherent randomness in MT-Bench evaluations.

Meanwhile, the performance of SimPO-multi lags behind that of SimPO, indicating that SimPO may not benefit from the multiple datasets. One possible reason is the lack of a reference model during training, which increases the model’s susceptibility to the inherent dispersion in the distribution of multiple datasets, ultimately leading to a less stable learned policy. Overall, these improvements suggest that by leveraging multiple datasets, UAPO-multi and SimUAPO-multi are able to learn a more accurate reference anchor that can extract deeper information of the prompt, for classifying responses as either winning or losing. The enhanced anchor allows the policy model to better capture and understand the preferences of the data.

Notably, UAPO achieves the best overall performance while also reducing training time. Unlike DPO and SimPO, which require three separate backward passes for three data pairs, UAPO uses a single pass, making it significantly more efficient.

3.3 In-depth Analysis of UAPO

The utility anchor effectively classifies responses of either winning or losing that closely relate to the prompt As one of the core contributions of this work, we introduce the utility anchor to of-

fine preference optimization. To demonstrate the classification performance of the utility anchor, we present the training curve of SimUAPO on Gemma-2-Instruct in Figure 1. Figure 1(a) illustrates a visualization that the utility anchor lies between the reward of the winning response and the reward of the losing response, indicating that the anchor effectively classifies responses as either winning or losing. Furthermore, the utility anchor ensures that the reward for the winning response decreases at a slower rate compared to the reward for the losing response. As a result, the reward margin, shown in Figure 1(b), increases during the training process. Additionally, the training accuracy, depicted in Figure 1(c), increases rapidly, as one of the anchor’s effects is to enhance the accuracy of training

Since the utility anchor is based on the prompt, we also sample two intuitive examples to measure the degree of relativity between the prompt and the utility anchor. As shown in Table 4, we select two prompts and present the rewards corresponding to the winning and losing responses for both the prompt and the utility anchor. The results reveal that the reward for the winning response of Prompt B is lower than the reference anchor for Prompt A. This underscores the importance of setting a utility anchor based on the prompt, as a static anchor would lead to incorrect classification.

Superior Generalization of UAPO under OOD Settings on RewardBench and RewardBench 2 As shown in Table 5, UAPO and SimUAPO

Method	RewardBench					RewardBench 2						
	Score	Chat	Chat Hard	Safety	Reasoning	Score	Factuality	Precise IF	Math	Safety	Focus	Ties
Gemma-2-9B-Instruct												
DPO	<u>81.0</u>	93.6	63.2	83.1	84.0	43.1	40.4	30.6	42.8	55.1	69.2	20.7
SimPO	79.2	93.3	62.3	80.1	81.2	41.7	41.5	25.8	42.6	50.2	65.1	24.8
UAPO	81.3	93.8	65.4	84.4	81.5	45.4	41.7	29.4	43.0	53.2	68.9	36.1
SimUAPO	80.8	93.6	61.3	86.1	82.0	<u>44.0</u>	42.2	25.6	43.3	62.1	66.3	24.2
Llama-3-8B-Instruct												
DPO	70.0	93.3	41.7	75.5	68.9	32.2	26.9	30.9	39.3	48.2	28.2	19.6
SimPO	<u>72.5</u>	93.0	53.1	74.2	69.8	<u>35.8</u>	32.3	30.0	43.2	50.9	36.3	22.1
UAPO	70.7	94.4	42.6	75.7	70.0	33.0	26.9	33.0	39.4	49.9	29.3	19.4
SimUAPO	72.7	93.9	53.2	75.3	68.5	36.2	33.2	30.1	38.9	51.0	33.8	29.9

Table 5: Results on out-of-distribution reward modeling benchmarks.

consistently outperform their respective baselines, DPO and SimPO, across multiple evaluation dimensions. On RewardBench, our methods exhibit stronger overall alignment quality, with notable improvements in more challenging aspects such as reasoning and safety. On RewardBench 2, which introduces more challenging tasks such as factuality, mathematical reasoning, and precise instruction following, UAPO and SimUAPO maintain a clear performance advantage. This indicates enhanced adaptability to distribution shifts and greater robustness in safety-critical evaluations.

These results demonstrate that UAPO not only improves preference alignment on standard benchmarks but also generalizes effectively to unseen scenarios, underscoring their robustness and practical applicability.

UAPO exhibits superior robustness under preference distribution shifts between training and evaluation To assess the robustness of preference optimization under distributional shifts, which commonly occur when aligning models with real-world human preferences, we simulate an out-of-distribution (OOD) training setting by randomly swapping the winning and losing responses in 40% of the training samples. This synthetic noise introduces preference corruption that reflects potential misalignment between training distributions and true human values (§2.2). We fine-tune both Mistral-Instruct and Llama-3-Instruct on this OOD dataset, and present the results in Table 6. The performance of baseline methods such as SimPO and SimUAPO declines due to the presence of corrupted preference signals. In contrast, SimUAPO shows a smaller performance drop, suggesting that it is more robust to noisy preferences. We attribute this improvement to the utility anchor mechanism, which encourages the model to focus on the in-

herent semantics of the prompt rather than relying heavily on uncertain or inconsistent human feedback.

Method	Mistral-7B-Instruct		Llama-3-8B-Instruct	
	LC (%)	WR (%)	LC (%)	WR (%)
SimPO	26.4 (-9.0%)	27.2 (-14.7%)	46.9 (-6.4%)	43.9 (-2.9%)
SimUAPO	28.0 (-2.1%)	30.9 (-4.0%)	48.7 (-4.9%)	41.8 (-12.2%)

Table 6: Performance comparison of Mistral-Instruct and Llama-3-Instruct on AlpacaEval 2 under a distribution shift in preferences between training and evaluation data. Values in parentheses represent the relative performance degradation compared to models trained on the original training data.

3.4 KL divergence control in UAPO

During training, smooth KL divergence is essential for preserving the policy model’s original behavior. As shown in Figure 2(a) and Figure 2(b), UAPO and SimUAPO exhibit lower and more stable KL divergence compared to DPO and SimPO, especially at the beginning of training. This demonstrates UAPO’s and SimUAPO’s superior control over KL divergence. Overall, both methods outperform their counterparts in performance while maintaining smooth, low KL divergence throughout training.

4 Related Work

Aligning LLMs with human preferences is a critical area of research. As mentioned before, several lines of research have been proposed to address this challenge, broadly categorized into approaches with explicit reward models and implicit reward models.

Explicit reward model approaches. Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) is a method used

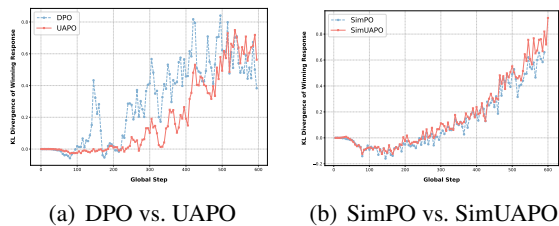


Figure 2: Analysis of KL divergence on Mistral-Instruct. (a) KL divergence of winning response on DPO and UAPO over training steps, and (b) KL divergence of winning response on SimPO and SimUAPO over training steps.

to fine-tune large language models (LLMs) by incorporating human-generated feedback. In RLHF, the process typically begins with training a reward model based on human preferences, where human evaluators provide feedback on model outputs. This reward model is then used to guide the reinforcement learning process, where the LLM’s parameters are adjusted to align its outputs with the human preferences. RLHF has been shown to significantly improve the alignment of LLMs, making their responses more aligned with human values and expectations. By iterating through this feedback loop, models can be fine-tuned to better match human decision-making and moral judgments. Building on this, Reinforcement Learning with AI Feedback (RLAIF) (Lee et al., 2023) extends the RLHF framework by replacing human feedback with AI-generated feedback. This innovative approach utilizes AI systems to evaluate and provide feedback on the model’s outputs, enabling alignment without the need for direct human input. RLAIF addresses some of the scalability challenges faced by RLHF, particularly in situations where large-scale human evaluation is impractical or costly. AI-generated feedback can be automated, allowing for faster iterations and potentially large-scale improvements in model behavior.

Implicit reward model approaches. Direct Preference Optimization (DPO) (Rafailov et al., 2023) simplifies the alignment process by directly optimizing the model based on human preferences, eliminating the need for a separate reward model. Identity Preference Optimization (IPO) (Azar et al., 2024) tried to resolve the issue of overfitting adding a regularization item. Contrastive Preference Optimization (CPO) (Xu et al., 2024) focuses on contrastive learning to align models with human preferences. By contrasting preferred outputs against less

preferred ones, CPO effectively guides the model toward desired behaviors. Kahneman-Tversky Optimization (KTO) (Ethayarajh et al., 2024) relies simply on binary feedback, which are more readily obtainable and more scalable for large datasets, as it reduces the need for pairwise data. Odds Ratio Preference Optimization (ORPO) (Hong et al., 2024) removes the need for a reference model and integrates supervised fine-tuning (SFT) with alignment into a single step. This approach simplifies the training process while maintaining alignment performance. R-DPO (Park et al., 2024) extends DPO by incorporating length normalization to disentangle the effects of response length and quality. This method helps improve the robustness of DPO in scenarios where response length may influence preference judgments. Simple Preference Optimization (SimPO) (Meng et al., 2024) is a variant of DPO that uses a reference-free reward with response length normalization. SimPO has demonstrated superior performance across diverse tasks by introducing a target reward margin to ensure a distinction between winning and losing responses. Robust DPO (rDPO) (Chowdhury et al., 2024) proposes an unbiased loss function to mitigate the impact of known label noise in pairwise preferences, while Distributionally Robustifying DPO (Dr.DPO) (Wu et al., 2025) adopts distributionally robust optimization to reweight samples and address both pointwise and pairwise noise effectively.

5 Conclusion

In this paper, we propose a general framework, UAPO, for offline preference optimization, which introduces a novel component, the uncertainty-aware utility anchor, to enhance learning stability and effectiveness. We validate UAPO on both DPO and SimPO, demonstrating that it achieves strong and consistent performance across three benchmarks. Unlike traditional approaches that rely solely on pairwise preference data, UAPO can incorporate multiple pair data, leading to significant improvement. These findings highlight the versatility and robustness of UAPO, making it a promising approach for real-world preference optimization tasks, especially in settings with limited data or computational resources.

Limitations

While UAPO serves as a general framework for offline preference optimization, we have only ver-

ified its performance on DPO and SimPO. Additionally, since UAPO introduces a utility anchor to preference optimization, it is crucial to investigate whether this anchor can further enhance the model through self-improvement. Future work could explore its application to other methods within this framework to further validate its generality.

Acknowledgments

This work was done during an internship at State Key Laboratory of General Artificial Intelligence, BIGAI. And this work was supported by the National Natural Science Foundation of China (62376031, 62337001), the Key Technologies R & D Program of Anhui Province (No. 202423k09020039), and the Fundamental Research Funds for the Central Universities.

References

AI@Meta. 2024. [Llama 3 model card](#).

Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR.

Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.

Sayak Ray Chowdhury, Anush Kini, and Nagarajan Natarajan. 2024. [Provably robust DPO: aligning language models with noisy feedback](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *CoRR*, abs/2110.14168.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. [Ultrafeedback: Boosting language models with high-quality feedback](#). *Preprint*, arXiv:2310.01377.

Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.

Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. [Kto: Model alignment as prospect theoretic optimization](#). *arXiv preprint arXiv:2402.01306*.

Adam Fisch, Jacob Eisenstein, Vicky Zayats, Alekh Agarwal, Ahmad Beirami, Chirag Nagpal, Pete Shaw, and Jonathan Berant. 2024. [Robust preference optimization through reward model distillation](#). *arXiv preprint arXiv:2405.19316*.

Leo Gao, John Schulman, and Jacob Hilton. 2023. [Scaling laws for reward model overoptimization](#). In *International Conference on Machine Learning*, pages 10835–10866. PMLR.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. [The language model evaluation harness](#).

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*.

Jiwoo Hong, Noah Lee, and James Thorne. 2024. [Orpo: Monolithic preference optimization without reference model](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11170–11189.

Jingru Jia, Zehua Yuan, Junhao Pan, Paul E McNamara, and Deming Chen. 2024. [Decision-making behavior evaluation framework for llms under uncertain context](#). *arXiv preprint arXiv:2406.05972*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. [Mistral 7b](#). *arXiv preprint arXiv:2310.06825*.

Ying Jin, Zhuoran Yang, and Zhaoran Wang. 2021. [Is pessimism provably efficient for offline rl?](#) In *International Conference on Machine Learning*, pages 5084–5096. PMLR.

Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Nathan Lambert, Valentina Pyatkin, Jacob Morrison, Lester James V. Miranda, Bill Yuchen Lin, Khyathi Raghavi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. [Rewardbench: Evaluating reward models for language modeling](#). In *Findings of the*

- Association for Computational Linguistics: NAACL 2025, Albuquerque, New Mexico, USA, April 29 - May 4, 2025*, pages 1755–1797. Association for Computational Linguistics.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Ren Lu, and Thomas Mesnard. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2310.07347*.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. 2024a. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. 2024b. [From live data to high-quality benchmarks: The arena-hard pipeline](#).
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 3214–3252. Association for Computational Linguistics.
- Saumya Malik, Valentina Pyatkin, Sander Land, Jacob Morrison, Noah A. Smith, Hannaneh Hajishirzi, and Nathan Lambert. 2025. [Rewardbench 2: Advancing reward model evaluation](#). *CoRR*, abs/2506.01937.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.
- Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. 2024. [Disentangling length from quality in direct preference optimization](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 4998–5017. Association for Computational Linguistics.
- Jan Peters and Stefan Schaal. 2007. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pages 745–750.
- Robin L Plackett. 1975. The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 24(2):193–202.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, and Christopher D Manning. 2023. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.16681*.
- Herbert A Simon. 1955. A behavioral model of rational choice. *The quarterly journal of economics*, pages 99–118.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshov, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kupala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan,

- Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024. [Gemma 2: Improving open language models at a practical size.](#)
- Amos Tversky and Daniel Kahneman. 1992. [Advances in prospect theory: Cumulative representation of uncertainty.](#) *Journal of Risk and Uncertainty*, 5(4):297–323.
- Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jiawei Chen, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. 2025. [Towards robust alignment of language models: Distributionally robustifying direct preference optimization.](#) In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. [Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation.](#) *arXiv preprint arXiv:2401.08417*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [Hellaswag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 4791–4800. Association for Computational Linguistics.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena.](#) In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

A Derivation of Eq. (11)

In Eq. (9), the sigmoid function $\sigma(\cdot)$ can be expanded, and the loss can be rewritten as:

$$\mathcal{L}_{\text{UAPO}}(\pi_\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \frac{\exp(r(x, y_w))}{\exp(r(x, y_\perp)) + \exp(r(x, y_w))} + \log \frac{\exp(r(x, y_\perp))}{\exp(r(x, y_\perp)) + \exp(r(x, y_l))} \right].$$

This expression can be decomposed into two terms:

$$\mathcal{L}_w = -\log \left(\frac{\exp(r(x, y_w))}{\exp(r(x, y_\perp)) + \exp(r(x, y_w))} \right), \quad \mathcal{L}_l = -\log \left(\frac{\exp(r(x, y_\perp))}{\exp(r(x, y_\perp)) + \exp(r(x, y_l))} \right).$$

When a given prompt is associated with multiple winning responses ($y_{w1}, y_{w2}, \dots, y_{wn}$) and multiple losing responses ($y_{l1}, y_{l2}, \dots, y_{lm}$), the reward of each winning response should exceed the utility anchor, which in turn should exceed the reward of each losing response. Under this setting, the two terms can be generalized as:

$$\mathcal{L}_w = -\sum_{y_w} \log \left(\frac{\exp(r(x, y_w))}{\exp(r(x, y_\perp)) + \sum_{y'_w} \exp(r(x, y'_w))} \right),$$

$$\mathcal{L}_l = -\log \left(\frac{\exp(r(x, y_\perp))}{\exp(r(x, y_\perp)) + \sum_{y_l} \exp(r(x, y_l))} \right).$$

Therefore, the UAPO objective with unpaired data can be expressed as Eq. (11).

B Derivation of Eq. (15)

A KL constrained policy optimization has been substantially defined by prior works (Rafailov et al., 2023). Following prior works (Jin et al., 2021), we can define a pessimistic function with reward uncertainty as a lower confidence bound. Typically, it has a form:

$$\max_{\theta} \mathbb{E}_{\mathcal{D}} [r(x) - \hat{u}(x)] - \beta \mathbb{D}_{\text{KL}}[\pi_\theta(x) \parallel \pi_{\text{ref}}(x)],$$

where $\hat{u}(x)$ is a uncertainty penalisation for input x .

Following prior work (Peters and Schaal, 2007), an optimal solution to the KL-constrained reward maximization:

$$\pi(y|x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left(\frac{1}{\beta} (r(x, y) - \hat{u}(x)) \right).$$

Taking the above optimal function into the Bradley-Terry from as in DPO (Rafailov et al., 2023), we have:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta) = -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} - \hat{u}(x) + \hat{u}(x) \right) \right].$$

Taking the implicit preference objective as in Rafailov et al. (2023) of pair y_w, y_l , we have the following inequalities:

$$\hat{r}(x, y_w) = \beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \geq \hat{u}(x) \geq \hat{r}(x, y_l) = \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)}.$$

Since the sigmoid function is monotonically increasing, we have

$$\log \sigma(\hat{r}(x, y_w) - \hat{u}(x)) + \log \sigma(\hat{u}(x) - \hat{r}(x, y_l)) \leq \log \sigma(\hat{r}(x, y_w) - \hat{r}(x, y_l)),$$

$$\mathcal{L}_{\text{DPO}}(\pi_\theta) \leq -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \hat{u}(x) \right) + \log \sigma \left(\hat{u}(x) - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right].$$

Therefore, UAPO in Eq. (9) can be considered as a lower bound optimization of the above uncertainty-aware optimization.

C Practical Implementation of Utility Anchors

For UAPO, we set

$$\log \pi(x, y_{\perp}) = \log \prod_{i=1}^n \sigma(Wh_i + b).$$

For SimUAPO, we set

$$\log \pi(x, y_{\perp}) = \frac{1}{n} \log \prod_{i=1}^n \sigma(Wh_i + b),$$

where $W \in \mathbb{R}^{1 \times d}$, $b \in \mathbb{R}$, σ is the sigmoid function, d is the hidden size of last layer, n denotes the number of tokens in the prompt and h_i denotes the output of the last hidden state corresponding to the i -th token of the prompt x .

For the representation of y_{\perp} , we assume it depends solely on the prompt x . Therefore, we use the hidden states of x as input to a linear function, whose output is the utility anchor. The weights of this linear function are learnable and updated during the preference optimization process, rather than being fixed in advance. This method allows the utility anchor to adjust according to the nuances of different prompts, thereby improving the alignment of the model’s outputs with human preferences.

In this paper, we use DPO and SimPO reward function to verify the feasibility of UAPO, *i.e.* replace $r(x, y)$ in Eq. (9) and Eq. (11) with $r_{\text{DPO}}(x, y)$ or $r_{\text{SimPO}}(x, y)$, denoted as UAPO and SimUAPO.

D Multiple Data Construction

Following (Meng et al., 2024), we construct the multiple dataset using data from princeton-nlp/mistral-instruct-ultrafeedback, princeton-nlp/llama3-ultrafeedback-armorm, and princeton-nlp/gemma2-ultrafeedback-armorm. In these datasets, for each prompt x , five responses are generated by the SFT model with a sampling temperature of 0.8. Next, the reward model annotates these responses to identify the highest-scoring one, denoted as y_w , and the remaining four responses are considered as y_l . To construct the multiple dataset, we select the y_w response and randomly choose three out of the four y_l responses. An example of this process can be seen in Table 7.

Prompt	how can i develop a habit of drawing daily
The Winning Response	<p>Developing a daily habit of drawing takes time and dedication, but it's a rewarding journey. Here's a breakdown of how to make it happen:</p> <ol style="list-style-type: none"> 1. Start Small and Be Realistic: <u>Tiny Habits:</u> Don't aim for masterpieces right away. Begin with 5-10 minutes a day. You'd be surprised how much you can accomplish even in short bursts. <u>Easy Wins:</u> Choose simple subjects to start with: doodles, shapes, still lifes, or even tracing. Build confidence with achievable goals. 2. Schedule It In: <u>Treat it like an appointment:</u> Choose a specific time each day dedicated to drawing. Put it in your calendar and stick to it as much as possible. <u>Morning Ritual:</u> Drawing first thing in the morning can set a creative tone for the day. <u>Evening Wind-Down:</u> If you're more relaxed in the evenings, drawing can be a calming way to de-stress. 3. Create a Conducive Environment: <u>Designated Space:</u> Set up a comfortable and inspiring drawing area. <u>Gather Supplies:</u> Have your pencils, paper, erasers, etc., readily available. 4. Remove Distractions: <u>Phone on Silent:</u> Put your phone away during your drawing time. <u>Quiet Space:</u> Find a spot where you won't be interrupted. 5. Embrace Imperfection: <u>Don't Edit:</u> Let go of the need for perfection. Focus on the process, not the outcome. <u>Experiment:</u> Try different styles, tools, and subjects. Don't be afraid to make mistakes—they're part of the learning process. 6. Find Inspiration: <u>Follow Artists:</u> Look at the work of artists you admire. <u>Take Photos:</u> Capture interesting scenes or objects to draw later. <u>Visit Museums:</u> Observe art in person. 7. Join a Community: <u>Online Groups:</u> Connect with other artists for support, feedback, and inspiration. <u>Local Classes:</u> Consider taking a drawing class to learn new techniques and meet fellow artists. 8. Track Your Progress: <u>Sketchbook:</u> Keep a sketchbook to record your daily drawings. Seeing your improvement over time can be incredibly motivating. <u>Photo Journal:</u> Take photos of your finished drawings to document your journey. 9. Reward Yourself: <u>Celebrate Milestones:</u> Treat yourself when you reach a drawing goal. <u>Enjoy the Process:</u> The most important reward is the joy of creating! <p>Remember, consistency is key. Even if you only draw for a few minutes each day, you'll see progress over time. Be patient, be persistent, and most importantly, have fun!</p>

**The First
Losing Response**

Developing a daily drawing habit takes dedication and a sprinkle of strategy. Here's a roadmap to help you on your artistic journey:

1. Set Realistic Goals: Start Small: Don't pressure yourself to create masterpieces every day. Begin with 15-20 minutes of drawing. Focus on Consistency: Aim for daily practice, even if it's just quick sketches. It's about building the habit, not producing perfect artwork.

2. Make it Enjoyable: Choose Subjects You Love: Draw things that genuinely interest you – people, animals, landscapes, objects, patterns – anything that sparks your creativity. Experiment with Mediums: Don't limit yourself to pencils. Explore pens, markers, watercolors, digital tools, or even charcoal. Listen to Music or Podcasts: Create a soundtrack or listen to something engaging while you draw to keep your mind stimulated.

3. Create a Routine: Schedule It In: Treat your drawing time like an important appointment. Block it out on your calendar and stick to it as much as possible. Find Your Ideal Time: Experiment with different times of day to see when you're most creative. Set Up a Dedicated Space: Even a small corner of your room can become your art haven. Keep your drawing supplies organized and easily accessible.

4. Embrace the Process: Don't Be Afraid to Make Mistakes: Mistakes are part of learning. Embrace them as opportunities to grow. Focus on Progress, Not Perfection: Celebrate small victories and track your improvement over time. Learn from Others: Study the work of artists you admire, take online courses, or join a drawing group for inspiration and feedback.

5. Stay Motivated: Join a Challenge: There are many online drawing challenges that can provide structure and accountability. Share Your Work: Show your drawings to friends, family, or online communities. Feedback and encouragement can be incredibly motivating. Reward Yourself: Celebrate milestones with small rewards to keep yourself on track.

Remember: Developing a daily habit takes time and effort. Be patient with yourself, stay consistent, and most importantly, have fun!

The Second Losing Response

It's awesome that you want to make drawing a daily habit! It's a rewarding journey. Here's a breakdown of how to make it stick:

1. Start Small & Be Realistic: Don't aim for masterpieces: Begin with 5-10 minutes a day. Even doodling is progress! Set achievable goals: Instead of "draw for an hour," try "sketch something for 15 minutes." Be kind to yourself: Some days will be better than others. Don't get discouraged if you miss a day or two.

2. Make it Enjoyable: Choose subjects you love: Draw things that excite you, whether it's animals, landscapes, people, or abstract shapes. Experiment with different mediums: Explore pencils, pens, charcoal, watercolors, digital tools – find what feels right. Listen to music or podcasts: Create a relaxing or inspiring atmosphere while you draw.

3. Build a Routine: Pick a specific time: Morning, lunch break, evening – find a time that works best for you and stick to it. Set reminders: Use your phone, calendar, or sticky notes to keep yourself on track. Create a dedicated space: Having a designated drawing area can signal to your brain that it's time to create.

4. Stay Inspired: Carry a sketchbook: Jot down ideas, quick sketches, or observations throughout the day. Follow artists you admire: Look at their work, learn from their techniques, and get inspired. Take online classes or workshops: Learn new skills and challenge yourself.

5. Embrace the Process: Focus on the journey, not just the destination: Enjoy the act of creating, even if the final result isn't perfect. Don't be afraid to make mistakes: They're opportunities to learn and grow. Celebrate your progress: Track your drawings, look back on how far you've come, and be proud of your dedication.

Remember, the key is to be consistent and patient with yourself. Enjoy the process of learning and growing as an artist!

The Third Losing Response

How to Develop a Daily Drawing Habit

It's fantastic that you want to make drawing a daily practice! Here's a roadmap to help you build that habit:

1. Start Small and Be Realistic: Don't aim for masterpieces: Begin with 5-10 minutes a day. Focus on consistency over quantity or quality initially. Set achievable goals: Instead of "draw for an hour," try "draw for 15 minutes before breakfast."

2. Make it Enjoyable: Choose subjects you love: Draw things that spark your interest, whether it's nature, people, animals, or abstract shapes. Experiment with different mediums: Don't limit yourself to pencils. Try pens, charcoal, watercolors, digital tools – explore and find what you enjoy most. Listen to music or podcasts: Find something that helps you relax and focus while you draw.

3. Create a Routine: Schedule it in: Treat drawing like any other important appointment. Block out a specific time in your day and stick to it. Set up your space: Have your drawing materials easily accessible in a dedicated area. Use a timer: Knowing you only have a set amount of time can help you get started and stay focused.

4. Embrace Imperfection: Don't be afraid to make mistakes: Mistakes are part of the learning process. Embrace them and learn from them. Don't compare yourself to others: Everyone progresses at their own pace. Focus on your own journey and celebrate your progress.

5. Find Inspiration: Follow artists you admire: Observe their techniques and styles. Visit museums or galleries: Expose yourself to different art forms and perspectives. Draw from life: Observe the world around you and capture what you see.

6. Track Your Progress and Celebrate Milestones: Keep a sketchbook: Document your daily drawings and track your progress over time. Share your work: Show your drawings to friends, family, or online communities for feedback and encouragement. Reward yourself: Acknowledge your achievements and celebrate your dedication.

Remember: Building a daily habit takes time and effort. Be patient with yourself, stay consistent, and enjoy the process!

Let me know if you have any other questions.

E Experiments

E.1 The robustness of UAPO over different constant reward anchors

We also evaluate the performance of SimUAPO and SimPO under different constant reward anchors (γ) settings on the Gemma-2-Instruct model in Table 8⁴. Across various constant reward anchors, SimUAPO consistently maintains high performance on AlpacaEval 2, demonstrating its robustness to changes in hyperparameters. In contrast, SimPO only achieves satisfactory performance when $\gamma = 5$, which fails to generate coherent sequences under other constant reward anchor settings.

⁴ denotes that the model has collapsed after training.

constant reward anchor γ	SimUAPO		SimPO	
	LC (%)	WR (%)	LC (%)	WR (%)
5	67.7	62.3	72.4	65.9
10	70.7	65.1	-	-
16	73.5	67.0	-	-
20	72.2	65.0	-	-

Table 8: Performance of using different constant reward anchors on Gemma-2-Instruct.

This performance disparity highlights the robustness of SimUAPO compared to SimPO. The key factor that contributes more is the utility anchor in SimUAPO, which effectively regulates the KL divergence between the policy model and the reference model. By dynamically adjusting the utility anchor, it is more robust to SimUAPO varying with hyperparameters changes (such as γ). This adaptability prevents excessive divergence from the reference model, which is more consistent and reliable performance across different training configurations.

E.2 The statistical significance testing of UAPO

To evaluate statistical significance, we conduct three evaluation runs on Mistral-Instruct using AlpacaEval 2 for SimPO, SimUAPO, and SimUAPO-multi, as shown in Table 9. The results show that UAPO is stable, as evidenced by its low standard deviation and notably low coefficient of variation, and that SimUAPO-multi consistently outperforms SimPO by a significant margin in terms of win rate.

Method	AlpacaEval 2 LC (%)	CV	AlpacaEval 2 WR (%)	CV
SimPO	28.868 \pm 0.147	0.509%	31.837 \pm 0.131	0.411%
SimUAPO	28.870 \pm 0.223	0.774%	32.257 \pm 0.021	0.065%
SimUAPO-multi	22.367 \pm 0.610	2.727%	40.463 \pm 0.254	0.628%

Table 9: The statistical significance testing results on AlpacaEval 2 with corresponding coefficient of variation (CV) on Mistral-Instruct.

E.3 Downstream tasks evaluation

We evaluate Mistral-Instruct and Llama-3-Instruct models trained with different methods on a range of tasks using the lm-evaluation-harness⁵ (Gao et al., 2024). These tasks include MMLU (Hendrycks et al., 2021), TruthfulQA (Lin et al., 2022), HellaSwag (Zellers et al., 2019), and GSM8K (Cobbe et al., 2021) and the results are shown in Table 10.

Knowledge is effectively retained. Both UAPO and SimUAPO preserve general knowledge as measured by MMLU, with minimal degradation compared to the SFT baseline. This suggests that the introduction of utility-aware preference objectives does not compromise the model’s foundational knowledge.

Truthfulness improves significantly. On the TruthfulQA benchmark, SimUAPO achieves the highest scores among all methods on both Mistral-Instruct and Llama-3-Instruct models, indicating that its preference optimization process emphasizes factual consistency and helps generate more truthful responses.

Commonsense reasoning remains competitive. For HellaSwag, both UAPO and SimUAPO maintain performance on par with or slightly better than SFT. This suggests that utility-aware optimization does not hinder the model’s ability to perform in-context reasoning and may benefit from exposure to diverse prompts in the preference dataset.

Math performance is more stable. Unlike many other methods that show sharp performance drops on GSM8K, UAPO demonstrates relatively stable results, suggesting that utility-aligned optimization may help balance the trade-off between human preference alignment and mathematical reasoning.

Overall performance is robust. UAPO consistently delivers competitive or improved average performance across benchmarks. And it balances strong gains in truthfulness and reasoning with only modest compromises in other areas, making it a strong candidate for robust preference alignment.

⁵<https://github.com/EleutherAI/lm-evaluation-harness>

Method	Mistral-7B-Instruct					Llama-3-8B-Instruct				
	MMLU(0)	TruthfulQA(0)	HellaSwag(0)	GSM8K(5)	Average	MMLU(0)	TruthfulQA(0)	HellaSwag(0)	GSM8K(5)	Average
SFT	59.06	66.87	83.67	42.23	62.96	63.79	51.56	75.88	75.28	66.63
DPO	59.27	66.74	84.53	42.76	63.33	64.61	55.50	76.23	76.42	68.19
SimPO	56.62	71.47	85.32	22.06	58.87	64.40	64.09	75.29	56.56	65.09
UAPO	57.78	65.45	83.39	30.17	59.20	64.31	55.25	76.34	76.04	68.00
SimUAPO	57.88	70.04	84.22	30.63	60.69	64.32	62.87	71.89	71.57	67.66

Table 10: Downstream task evaluation results of tasks on the huggingface open leaderboard.

Judgement	Score
The assistant’s response to the second question contains an error. It suggests that a solution ... The corrected rating for the assistant’s response would be [[5]].	5
The assistant’s response is inaccurate. Even though the assistant’s code meets the user’s requirement of ... Therefore, I would rate this response as: [[4]].	4

Table 11: Inconsistencies in MT-Bench scores for the same model on a single sample.

E.4 The example which shows the randomness of MT-Bench

During our evaluation process, we frequently observe that the same model can receive varying scores on a single MT-Bench sample (e.g., see Table 11). This variability highlights concerns regarding the consistency and reliability of MT-Bench. Consequently, we consider AlpacaEval 2 and Arena-Hard to offer more stable and trustworthy evaluations for comparing different preference optimization methods. Nevertheless, since MT-Bench remains widely used in prior work, we include it in our evaluation to ensure consistency and comparability across studies.

E.5 Multiple datasets performance on Mistral-Instruct

We also train DPO, SimPO, UAPO, and SimUAPO on Mistral-Instruct using multiple datasets and the results are in Table 12. SimUAPO-multi achieves a 40.5 WR on AlpacaEval 2 and a 23.1 WR on Arena-Hard which are much higher than other preference optimization methods. Interestingly, on Mistral-Instruct, SimUAPO-multi further enlarges the gap between LC and WR on AlpacaEval 2 compared to both SimPO and SimUAPO. This discrepancy is not observed on Llama-3-Instruct, which may be attributed to the inherent characteristics of the Mistral-Instruct model and the specific design of the SimPO reward function.

Method	AlpacaEval 2		Arena-Hard	MT-Bench
	LC (%)	WR (%)	WR (%)	GPT-4
multiple datasets				
DPO	20.5	18.1	13.4	<u>7.6</u>
SimPO	29.0	31.9	<u>19.9</u>	<u>7.2</u>
UAPO	23.1	17.7	11.1	7.7
SimUAPO	<u>28.6</u>	<u>32.2</u>	17.7	<u>7.6</u>
DPO-multi	17.0	16.7	8.4	6.8
SimPO-multi	14.5	22.2	9.2	6.3
UAPO-multi	24.1	18.6	13.8	7.7
SimUAPO-multi	21.9	40.5	23.1	<u>7.6</u>

Table 12: Performance comparison of different methods on Mistral-Instruct trained with multiple datasets.

E.6 The necessity of constant reward anchor in UAPO

To further investigate the significance of the constant reward anchor in UAPO, we compare the performance of UAPO and SimUAPO with and without the constant reward anchor γ in Table 13. The results show that removing the constant reward anchor leads to a significant performance decline for both UAPO and SimUAPO across most benchmarks, underscoring its critical role in the UAPO framework. However, UAPO exhibits inconsistent trends without the reward anchor: while performance drops on AlpacaEval 2, it improves on Arena-Hard. This inconsistency across benchmarks when evaluating the same model highlights the need for more reliable and robust evaluation benchmarks for preference optimization.

Method	Mistral-7B-Instruct				Gemma-2-9B-Instruct			
	AlpacaEval 2		Arena-Hard	MT-Bench	AlpacaEval 2		Arena-Hard	MT-Bench
	LC (%)	WR (%)	WR (%)	GPT-4	LC (%)	WR (%)	WR (%)	GPT-4
w constant reward anchor								
UAPO	23.1	17.7	11.1	<u>7.7</u>	<u>70.2</u>	67.7	<u>58.8</u>	8.9
SimUAPO	28.6	32.2	<u>17.7</u>	<u>7.6</u>	73.5	<u>67.0</u>	59.4	8.9
w/o constant reward anchor								
UAPO	17.6	14.8	14.9	<u>7.6</u>	55.4	48.7	41.4	8.4
SimUAPO	<u>24.8</u>	<u>26.9</u>	22.7	<u>7.6</u>	60.8	56.6	54.9	<u>8.5</u>

Table 13: An ablation study of the constant reward anchor-free setting on Mistral-Instruct and Gemma-2-Instruct.

E.7 Generalization capability of UAPO on R-DPO

To further demonstrate the generality of UAPO, we conduct an additional experiment by applying it to R-DPO on Llama-3-Instruct. As shown in Table 14, R-UAPO achieves a 2.2-point improvement over R-DPO on AlpacaEval 2 LC. These results highlight that UAPO exhibits strong generalization capabilities.

Method	AlpacaEval 2 LC (%)	AlpacaEval 2 WR (%)
R-DPO	48.0	45.8
R-UAPO	50.2	47.1

Table 14: The performance of R-DPO and R-UAPO on Llama-3-Instruct trained with the pairwise dataset.

F Implementation Details

F.1 Training Hyperparameters

Since the performance of most offline preference optimization methods is highly sensitive to training hyperparameters, such as learning rate, batch size, and others, we conducted a hyperparameter search, as specified in the respective papers, to ensure a fair comparison. The results are summarized in Table 15. We use a batch size of 64 for training all methods. The learning rate is selected from the range $[3e - 7, 5e - 7, 7e - 7, 1e - 6]$. All models are trained for one epoch using the Adam optimizer (Kingma and Ba, 2015).

Method	Objective	Hyperparameter
DPO (Rafailov et al., 2023)	$-\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right)$	$\beta \in [0.01, 0.05, 0.1]$
IPO (Azar et al., 2024)	$\left(\log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} - \frac{1}{2\tau} \right)^2$	$\tau \in [0.01, 0.1, 0.5, 1.0]$
CPO (Xu et al., 2024)	$-\log \sigma \left(\beta \log \pi_{\theta}(y_w x) - \beta \log \pi_{\theta}(y_l x) \right) - \lambda \log \pi_{\theta}(y_w x)$	$\lambda = 1.0, \beta \in [0.01, 0.05, 0.1]$
KTO (Ethayarajh et al., 2024)	$-\lambda_w \sigma \left(\beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - z_0 \right) + \lambda_l \sigma \left(z_0 - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} \right)$, where $z_0 = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\beta \text{KL}(\pi_{\theta}(y x) \parallel \pi_{\text{ref}}(y x))]$	$\lambda_l = \lambda_w = 1.0$ $\beta \in [0.01, 0.05, 0.1]$
ORPO (Hong et al., 2024)	$-\log p_{\theta}(y_w x) - \lambda \log \sigma \left(\frac{p_{\theta}(y_w x)}{1 - p_{\theta}(y_w x)} \right) - \log \left(\frac{p_{\theta}(y_l x)}{1 - p_{\theta}(y_l x)} \right)$, where $p_{\theta}(y x) = \exp \left(\frac{1}{ y } \log \pi_{\theta}(y x) \right)$	$\lambda \in [0.1, 0.5, 1.0, 2.0]$
R-DPO (Park et al., 2024)	$-\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)} + (\alpha y_w - \alpha y_l) \right)$	$\alpha \in [0.05, 0.1, 0.5, 1.0]$ $\beta \in [0.01, 0.05, 0.1]$
SimPO (Meng et al., 2024)	$-\log \sigma \left(\frac{\beta}{ y_w } \log \pi_{\theta}(y_w x) - \frac{\beta}{ y_l } \log \pi_{\theta}(y_l x) - \gamma \right)$	$\beta \in [2.0, 2.5, 10.0]$ $\gamma \in [0.3, 1.0, 1.6, 3.0, 5.0]$
UAPO	$-\log \sigma \left(r_{\text{UAPO}}(x, y_w) - r_{\text{UAPO}}(x, y_{\perp}) \right) - \log \sigma \left(r_{\text{UAPO}}(x, y_{\perp}) - r_{\text{UAPO}}(x, y_l) \right)$	$\beta \in [0.01, 0.05], \gamma \in [1.0, 4.5, 8.0]$
SimUAPO	$-\log \sigma \left(r_{\text{SimUAPO}}(x, y_w) - r_{\text{SimUAPO}}(x, y_{\perp}) \right) - \log \sigma \left(r_{\text{SimUAPO}}(x, y_{\perp}) - r_{\text{SimUAPO}}(x, y_l) \right)$	$\beta \in [2.5, 10.0], \gamma \in [4.5, 8.0]$

Table 15: The objectives and hyperparameters of different offline preference optimization methods.

F.2 Computing Resources

All the experiments can be conducted on eight Nvidia A100 80GB GPUs, 32GB memory, and a 128-core AMD CPU.

F.3 Training Time

Since UAPO introduces additional computational overhead for optimizing and inferring dummy tokens, we compare the training time of different methods. All experiments are conducted with PyTorch 2.3.0 on Ubuntu 18.04. The results show the training time required for each method in Table 16.

Method	Mistral-7B-Instruct	Llama-3-8B-Instruct	Gemma-2-9B-Instruct
DPO	5.3h	5.0h	4.7h
SimPO	4.7h	5.3h	10.9h
UAPO	5.9h	5.3h	5.1h
SimUAPO	3.3h	6.0h	12.1h

Table 16: Training time comparison across different methods.