# When Human Preferences Flip: An Instance-Dependent Robust Loss for RLHF

**Yifan Xu,**[1] **Xichen Ye,**[2] **Yifan Chen,**[1,*] **Qiaosheng Zhang**[3]

[1]Hong Kong Baptist University
[2]Fudan University
[3]Shanghai Artificial Intelligence Laboratory

## Abstract

Quality of datasets plays an important role in large language model (LLM) alignment. In collecting human feedback, however, *preference flipping* is ubiquitous and causes corruption in data annotation; the issue necessitates the alignment algorithms with improved robustness against potential flipped pairs. To this end, this paper introduces a Flipping-Aware Direct Preference Optimization (FA-DPO) algorithm tailored to preference flipping from a reinforcement learning with human feedback (RLHF) perspective. We dissect the inherent human intention model and the preference flipping mechanism introduced by external factors as two distinct stages; in the latter, we introduce an instance-dependent flipping probability on the basis of the Bradley-Terry (BT) model. Further, by leveraging features relevant to preference annotation, we capture uncertainty in judgments and model preference flipping patterns. In practice, we design a simple yet efficient iterative optimization algorithm compatible with the original RLHF and DPO algorithms. In our experiments, we investigate the instance-dependent preference flipping model under multiple circumstances for evaluation of our proposed method, as well as other baseline methods.

## 1 Introduction

Alignment has been identified as a crucial approach contributing to the remarkable capacity of large language models (LLMs) to understand human intentions. In general, this approach enables LLMs to produce responses that align well with human expectations and reduce toxic generations (Achiam et al. 2023; Dubey et al. 2024). Among existing alignment paradigms, a notable one is the reinforcement learning from human feedback (RLHF), which, along with its numerous variants, has attracted significant attention (Ouyang et al. 2022; Casper et al. 2023).

Despite the effectiveness of current RLHF methods, they implicitly suffer from noise in human feedback data (Zheng et al. 2023; Gao, Alon, and Metzler 2024). For example, Gao, Alon, and Metzler (2024) reported that a $10\%$ increase in preference flipping ratios can result in a $30\%$ decrease in alignment performance, as measured by *the win rate* (defined in Section 5.1). As noise is inevitably intro-

duced during data collection process, robust alignment algorithms bring the benefits of not only defending potential dataset attacks, but also reducing the costs of collecting clean data. One existing genre of robust learning approaches leverages noise information inferred from existing data (Xia et al. 2021; Song et al. 2022), and this characteristic is especially attractve in the contexts of direct alignment algorithms such as direct preference optimization (DPO; Rafailov et al. 2024), where training relies on a fixed offline dataset.

In tackling robustness challenges in RLHF, previous research focuses on a simplified scenario where preferences are randomly flipped at a *fixed* rate (Chowdhury, Kini, and Natarajan 2024; Wu et al. 2024a; Cheng et al. 2024; Liang et al. 2024). However, in practical applications, it is ungrounded to assume the possibility of annotation errors is independent of the specific content being annotated. In this work, we instead investigate a setting of "instance-dependence" (Xia et al. 2020; Liu, Cheng, and Zhang 2023), which poses greater challenges as the noise distribution can vary significantly across different samples. Therefore, in this work, we aim to answer the following question:

> *How can we properly model the instance-dependent preference flipping incurred during preference data annotation?*

In response, we propose a statistically consistent approach to align with corrupted human feedback, which we refer to as Flipping-Aware Direct Preference Optimization (FA-DPO). This approach ensures that, given accurate estimation of the flipping model, the learned policy will achieve consistency comparable to learning the policy under clean data. We first explicitly model the flipping probability on the basis of the Bradley-Terry model, to connect with the posterior probabilities of the observed label.

Specifically, the annotation process under our framework can be viewed as two sequential stages: ❶ labeling according to true human intention and ❷ instance-dependent label contamination, where the instance-dependent preference flipping occurs after the labeling stage ❶, representing a transition from the true label to the flipped label for each sample. In this way, given a corrupted human preference dataset, FA-DPO can post-train the LLM parameters robustly and produce optimal policies, as with the true labels, via correction for the original BT model loss with the posterior of the observed corrupted labels.

---

In particular, we further introduce an iterative optimization framework that jointly optimizes both the flipping estimation model and the LLM in post-training. To capture the transition between the original label and the observed label, we leverage a classification module on relevant features to construct the flipping probability within a preference pair. The main contributions of this paper are three-fold:

1. We investigate a challenging RLHF / DPO setting with instance-dependent preference flipping, and propose a novel probabilistic model to characterize the noise.

2. We address the instance-dependent estimation of preference flipping probabilities via a classification model that incorporates informative preference features studied in natural language processing.

3. We propose a simple yet efficient iterative optimization algorithm based on RLHF and DPO to capture the noise in the offline dataset during post-training. The whole training pipeline is validated in our experiments.

## 2 Related Works

We review the previous works on preference alignment and (generalized) robust RLHF in this section. Due to context limitation, we leave the complete review to Appendix B for the reader's convenience.

**Preference alignment** The most well-known approach for preference alignment is Reinforcement Learning from Human Feedback (Ziegler et al. 2019; Ouyang et al. 2022, RLHF), which involves training a *reward model* to capture human preferences and then guiding LLMs to *generate high-reward responses using reinforcement learning algorithms* such as Proximal Policy Optimization (PPO) (Schulman et al. 2017). However, in practice RL-based methods can be complex and unstable during training (Rafailov et al. 2024; Wu et al. 2024c; Yuan et al. 2023). As a result, recent research has focused on simpler and more stable alternatives to RLHF, namely, direct preference alignment (Rafailov et al. 2024; Zhao et al. 2023; Ethayarajh et al. 2024; Azar et al. 2024; Meng, Xia, and Chen 2024).

**RLHF against perturbations** Current approaches of robust RLHF can be categorized into three main types. ❶ Noise fitting (Bukharin et al. 2024) involves making assumptions on the noise in the data and incorporating this modeling into the reward learning process, which will be jointly optimized with the parameterized reward function. ❷ Sample selection (or sample re-weighting) methods (Cheng et al. 2024) leverage this phenomenon to identify clean samples based on the loss values observed during the training process. Some approaches address corruption in input data through ❸ robust loss design (Chowdhury, Kini, and Natarajan 2024; Liang et al. 2024), focusing on constructing loss functions that are more resistant to data noise; this technique is also known as label smoothing (Mitchell 2023).

## 3 Preliminaries

In this section, we introduce the basics of RLHF for LLM alignment in advance of the formal proposal of our methodology in Section 4.

### 3.1 RLHF for LLM alignment

Let an input (prompt) of an LLM be $x \in \mathcal{X}$ and the generated output (response) be $y \in \mathcal{Y}$. In the formulation for RLHF (Rafailov et al. 2024), the LLM is viewed as a policy $\pi_\theta(y \mid x)$ parameterized by $\theta$, which outputs an action $y$ (response) based on the state $x$ (prompt). Preference data is further collected and annotated by human labelers, denoted as $y_w \succ y_l \mid x$, where $y_w$ is the preferred response and $y_l$ is the dispreferred one in $(y_1, y_2)$ for the prompt $x$.

The pre-trained model first undergoes on round of supervised fine-tuning (SFT), resulting in a reference model $\pi_{\text{ref}}$. Then the Bradley-Terry (BT) model (Bradley and Terry 1952) is employed to relate the preference data $\{y_w^i \succ y_l^i\}$ to a reward model $r(x, y)$. The connection is formulated as:

$$p^*(y_w \succ y_l \mid x) = \sigma(r^*(x, y_w) - r^*(x, y_l)), \quad (1)$$

where $\sigma$ is the standard sigmoid function, and $r^*(\cdot)$ is the optimal reward model. Using the BT model and the maximum likelihood principle, the loss for learning the reward model is:

$$\mathcal{L}_R(\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l)) \right], \quad (2)$$

where $r_\phi(\cdot)$ is the reward model parameterized by $\phi$.

After training the reward model $r_\phi(\cdot)$ from Equation (2), reinforcement learning (RL) is applied to optimize the LLM $\pi_\theta$ with the reward signals provided by $r_\phi(\cdot)$. The optimization objective for $\pi_\theta$ is formulated as:

$$\mathcal{L}_\pi(\theta) = -\mathbb{E}_{x \sim \mathbb{P}_x, y \sim \pi_\theta(\cdot \mid x)} [r_\phi(x, y)] + \beta \cdot \mathbb{E}_{x \sim \mathbb{P}_x} [\mathbb{D}_{\text{KL}}(\pi_\theta(\cdot \mid x) \| \pi_{\text{ref}}(\cdot \mid x))], \quad (3)$$

where $\mathbb{P}_x$ represents the marginal distribution of the prompt $x$. The first term corresponds to reward maximization in standard RL optimization. The second term is a KL divergence that constrains the update of $\pi_\theta$ to not deviate from the reference model $\pi_{\text{ref}}$, $\beta$ is the coefficient that weights the KL divergence between $\pi_\theta$ and $\pi_{\text{ref}}$.

### 3.2 DPO as efficient RLHF

Direct preference optimization (DPO) is an algorithm proposed for practical efficiency (compared to original RLHF), that merges the two stages in RLHF, reward modeling and policy optimization, into a single step. In DPO, the policy is optimized directly using the offline dataset $\mathcal{D}$ without constructing an explicit reward model.

As shown by Peng et al. (2019), the closed-form solution of the conditional distribution $\pi(y \mid x)$ that minimizes Equation (3) is:

$$\pi_r(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r_\phi(x, y)\right), \quad (4)$$

where $Z(x) = \sum_y \pi_{\text{ref}}(y \mid x) \exp(1/\beta \cdot r_\phi(x, y))$ is the partition function for normalization. By applying this result to Equation (2), the DPO loss function that includes only the parameterized $\pi_\theta$ as the optimization variable is derived as:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma(\hat{r}_\theta(x, y_w) - \hat{r}_\theta(x, y_l)) \right], \quad (5)$$

where $\hat{r}_\theta(x, y) = \beta \log(\pi_\theta(y|x)/\pi_{\text{ref}}(y|x))$ is the implicit reward model derived from $\pi_\theta$.

# 4 Method

In this section, we begin by providing an overview of the motivation to model instance-dependent preference flipping. We then formulate the flipping procedure under the RLHF framework, and illustrate how this estimation can be integrated into the standard RLHF pipeline. Finally, we present a detailed model design for the instance-dependent preference flipping probability estimation in FA-DPO.

## 4.1 Motivations

It is widely accepted that in real-world scenarios, the intrinsic preference labeling mechanism of humans can be unified as a general human intention model, represented by the Bradley-Terry model (Bai et al. 2022). Noisy human feedback, however, usually stems from external factors rather than inherent errors in human decisions. For instance, environmental distractions may compromise an annotator's focus, reducing labeling accuracy, and external noise can also be introduced by maliciously altering original annotations.

Consequently, this corruption process can be viewed as a post-transition applied to the initial labeling mechanism governed by the human intention model. We model this process as instance-dependent preference flipping, positing that the probability of a flip is correlated with the data content as well as the relation between each pair of responses.

We dig into the statistical structure of the flipping pattern in the preference data as opposed to the noise sparsity assumption in (Bukharin et al. 2024). The proposed modeling, as stated in Section 4.2, is supposed to exploit the relation between the noisy preference posterior and the true likelihoods. The neural network module to capture the noise is introduced in Section 4.3.

## 4.2 Flipping-aware Loss

We start the derivation of the flipping-aware loss from the RLHF model, and then extend the result to the DPO setting. Given a corrupted dataset $\tilde{\mathcal{D}}$, for each noisy triplet $(x, \tilde{y}_w, \tilde{y}_l)$, we denote the sample as $\tilde{x}$ for simplicity. Following the principle of maximum likelihood, the usual loss for preference modeling is:

$$\mathcal{L}_{\text{MLE}} = -\mathbb{E}_{(x, \tilde{y}_w, \tilde{y}_l) \sim \tilde{\mathcal{D}}}[\log \mathbb{P}\{\tilde{y}_w \succ \tilde{y}_l \mid x\}]. \quad (6)$$

In standard RLHF, we only have the parameterized preference probability under clean data, i.e., $\mathbb{P}\{y_w \succ y_l \mid x\}$; the direct usage of the standard RLHF loss in Equation (6) can be sub-optimal. To ease the discussion of how to bridge the gap between preference distributions under clean and noisy data, we first formalize the preference flipping process through the following proposition:

**Proposition 4.1** (Instance-dependent preference flipping). *For any input $\tilde{x}$, the corrupted preference probability, under the instance-dependent preference flipping setting, relates to the true preference likelihood via:*

$$\tilde{\mathbb{P}}\{\tilde{y}_w \succ \tilde{y}_l \mid x\} = (1 - \varepsilon_{\tilde{x}})p + \varepsilon_{\tilde{x}}(1 - p),$$

*where $\varepsilon_{\tilde{x}}$ represents the instance-specific flipping probability for triplet $\tilde{x} = (x, \tilde{y}_w, \tilde{y}_l)$, and $p$ denotes the true likelihood $\mathbb{P}\{\tilde{y}_w \succ \tilde{y}_l \mid x\}$ for brevity.*

The proof of Theorem 4.1 is direct and omitted. We can then establish the relation between the corrupted posterior $\tilde{\mathbb{P}}\{\tilde{y}_w \succ \tilde{y}_l \mid x\}$ observed in noisy data and the underlying clean probability $\mathbb{P}\{\tilde{y}_w \succ \tilde{y}_l \mid x\}$. This enables us to recover the true preference probabilities by accounting for the instance-dependent flipping process:

$$\mathcal{L}_{\text{FA-DPO}} = -\mathbb{E}_{\tilde{x} \sim \tilde{\mathcal{D}}}\left[\log\left((1 - \varepsilon_{\tilde{x}})p + \varepsilon_{\tilde{x}}(1 - p)\right)\right], \quad (7)$$

which reads the loss design in FA-DPO. For practical implementation with the BT model in reward-based RLHF, we parameterize $p$ with $\phi$ and denote:

$$p_\phi = \sigma(r_\phi(x, \tilde{y}_w) - r_\phi(x, \tilde{y}_l)).$$

Similarly, for the DPO parameterization, the preference probability is given by substituting the probabilistic modeling Equation (4) into the formula above:

$$p_\theta = \sigma\left(\beta \log \frac{\pi_\theta(\tilde{y}_w \mid x)}{\pi_{\text{ref}}(\tilde{y}_w \mid x)} - \beta \log \frac{\pi_\theta(\tilde{y}_l \mid x)}{\pi_{\text{ref}}(\tilde{y}_l \mid x)}\right).$$

**Comparison with cDPO and rDPO.** To better understand FA-DPO, we compare it with related approaches, cDPO (Mitchell 2023) and rDPO (Chowdhury, Kini, and Natarajan 2024), from the perspective of gradients.

First, we recall standard DPO directly substitutes the paramterized preference probability $\mathbb{P}_\theta\{\tilde{y}_w \succ \tilde{y}_l \mid x\}$ for the corrupted posterior $\tilde{\mathbb{P}}\{\tilde{y}_w \succ \tilde{y}_l \mid x\}$, ignoring the preference flipping mechanism underlying the corrupted datasets. Other robust losses, such as cDPO (Mitchell 2023), instead similarly follow the principle of maximum likelihood and adopt a loss correction technique, but the $\varepsilon$ parameter in cDPO is a hyperparameter irrelevant to each sample (while our loss considers an instance-dependent preference flipping setting). On the basis of cDPO, rDPO (Chowdhury, Kini, and Natarajan 2024) further debias the loss function.

To better understand the difference between these methods, we compare the gradient weights of these methods with that of FA-DPO.

Following the notations in Chowdhury, Kini, and Natarajan (2024), we characterize the gradient of DPO-like methods as

$$\nabla_\theta \mathcal{L} = -\zeta \cdot \beta \left(\nabla_\theta \log \pi_\theta(\tilde{y}_w \mid x) - \nabla_\theta \log \pi_\theta(\tilde{y}_l \mid x)\right). \quad (8)$$

The weighting coefficient $\zeta$ for each method can be related to the base DPO weight coefficient $\zeta_{\text{DPO}} = 1 - p_\theta$, where cDPO applies a small reduction $\zeta_{\text{cDPO}} = \zeta_{\text{DPO}} - \varepsilon$ while rDPO introduces an additive correction $\zeta_{\text{rDPO}} = \zeta_{\text{DPO}} + \frac{\varepsilon}{1-2\varepsilon}$ to adjust the weighting behaviors, respectively. We specify the coefficient $\zeta_{\text{FA-DPO}}$ in FA-DPO as follows. The derivation is deferred in Appendix A.

**Lemma 4.2** (Gradient weight coefficient). *For a triplet $\tilde{x} = (x, \tilde{y}_w, \tilde{y}_l)$, we have the gradient weight coefficient for FA-DPO as*

$$\zeta_{\text{FA-DPO}} = \frac{(1 - 2\varepsilon_{\tilde{x}})p_\theta}{(1 - 2\varepsilon_{\tilde{x}})p_\theta + \varepsilon_{\tilde{x}}} \cdot \zeta_{\text{DPO}}.$$

The above lemma demonstrates that our weighting scheme constitutes a reparametrization of the DPO gradient weight, distinct from the additive correction approaches

employed by cDPO and rDPO. Crucially, when no flipping occurs ($\varepsilon_{\tilde{\mathbf{x}}} = 0$), the gradient reduces exactly to that of the standard DPO. For low flipping probabilities ($\varepsilon_{\tilde{\mathbf{x}}} < 0.5$), the weight increases with the model's confidence ($p_\theta$), enhancing stability during convergence and improving robustness to noise compared to the fixed correction mechanisms of cDPO and rDPO. At a flipping probability of 0.5, indicating inherent ambiguity in the preference signal, the weight becomes zero, automatically filtering out these low-margin samples that could otherwise impair model performance. Most significantly, when $\varepsilon_{\tilde{\mathbf{x}}} > 0.5$, the gradient direction reverses to optimize $\mathbb{P}\{\tilde{y}_l \succ \tilde{y}_w \mid x\}$ instead of $\mathbb{P}\{\tilde{y}_w \succ \tilde{y}_l \mid x\}$ under conditions of high model uncertainty. This demonstrates self-correction for samples with high detected flipping rates—a capability absent in cDPO and rDPO which apply uniform corrections. Furthermore, the weight approaches zero when high model confidence contradicts high estimated $\varepsilon_{\tilde{\mathbf{x}}}$, yielding robust weights that are jointly determined by both the flipping probability and the policy model's confidence.

### 4.3 Transition Probability Modeling

**Preference flipping modeling**  To overcome the limitations of existing robust DPO approaches that assume either (1) a uniform flipping ratio across all samples or (2) sparse noise patterns within a given dataset, we introduce an instance-dependent preference flipping module that dynamically estimates sample-specific flipping probabilities based on instance features. we model this probability as a logistic regression function of input-dependent features:

$$\varepsilon_{\tilde{\boldsymbol{x}}} = \sigma(\langle \omega, h(\tilde{\boldsymbol{x}}) \rangle + \omega_0), \tag{9}$$

where $h : \mathcal{X} \to \mathbb{R}^d$ is the feature map and $\omega \in \mathbb{R}^d$ are learnable parameters.

**Feature map construction**  We deliberately design the feature map $h(\cdot)$ to incorporate three concepts validated to be effective in language modeling. Notably, the features are supposed to be permutation-equivariant to the response pairs, since the order thereof can be arbitrary in the corrupted sample triplet $\tilde{\boldsymbol{x}} = (x, \tilde{y}_w, \tilde{y}_l)$.

*Response Length.* It is noted that longer responses increase cognitive load for human annotator, which may increase the error rate (Chen et al. 2024). We compute both the average lengths and the length difference within a preference sample triplet $(x, \tilde{y}_w, \tilde{y}_l)$ as the first feature,

$$h_{\text{len}}(\tilde{\boldsymbol{x}}) = \left[ \frac{|\tilde{y}_w| + |\tilde{y}_l|}{2}, ||\tilde{y}_w| - |\tilde{y}_l|| \right]^\top.$$

Here, even we change the order of the response pair to $(x, \tilde{y}_l, \tilde{y}_w)$, the feature map $h_{\text{len}}(\cdot)$ is invariant.

*Perplexity (PPL).* Perplexity reflects the uncertainty or complexity of a probability distribution. High PPL often correlates with sentences that are hard for humans to comprehend, thereby increasing annotation difficulty (Kong et al. 2024).

$$h_{\text{ppl}}(\tilde{\boldsymbol{x}}) = \left[ \frac{\log(\pi_\theta(\tilde{y}_w|x)\pi_\theta(\tilde{y}_l|x))}{2}, \left| \log \frac{\pi_\theta(\tilde{y}_w|x)}{\pi_\theta(\tilde{y}_l|x)} \right| \right]^\top.$$

*Reward Margin.* The reward margin implicitly quantifies the model's confidence in distinguishing preferred from dispreferred responses, which is often used for data selection in preference learning (Wu et al. 2024b; Huang et al. 2025). The corresponding feature map is

$$h_{\text{margin}}(\tilde{\boldsymbol{x}}) = \left[ \frac{\hat{r}_\theta(x, \tilde{y}_w) + \hat{r}_\theta(x, \tilde{y}_l)}{2}, |\hat{r}_\theta(x, \tilde{y}_w) - \hat{r}_\theta(x, \tilde{y}_l)| \right]^\top$$

where $\hat{r}_\theta(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$ is the implicit reward function induced by DPO.

These features are then concatenated, along with a scalar 1 (for the bias term in $\omega$), and scaled to form the final feature map:

$$h(\tilde{\boldsymbol{x}}) = [h_{\text{len}}(\tilde{\boldsymbol{x}}), h_{\text{ppl}}(\tilde{\boldsymbol{x}}), h_{\text{margin}}(\tilde{\boldsymbol{x}})]^\top \tag{10}$$

**Iterative update**  To optimize both the preference flipping model and the LLM, we design an iterative update paradigm, which is also adopted in previous work for robust RLHF (Bukharin et al. 2024).

Regarding the concrete design, empirical studies have demonstrated that deep neural networks exhibit a consistent learning trajectory, initially capturing generalizable patterns before eventually overfitting to noisy training instances (Cheng et al. 2024). To adapt to this characteristic, in practice we utilize the learned capability of the LLM during the initial stage of training to optimize the flipping model first (this operation is referred to as *warmup*), and then we iteratively update the two models. The complete training algorithm is presented in Algorithm 1 in Appendix D.

### 4.4 Theoretical Analysis

This section establishes the theoretical foundations of our approach. Specifically, we demonstrate that the proposed loss function for preference flipping yields desirable statistical properties: consistency under noise and convergence guarantees for the flipping model parameters.

We first present the key result that the minimizer of our FA-DPO loss function, operating on the corrupted (flipped preference) data distribution $\tilde{\mathcal{D}}$, coincides with the minimizer of the original loss function on the underlying clean data distribution $\mathcal{D}$ in RLHF. This *consistency* property guarantees that, asymptotically, our method recovers the same optimal model parameters as would be obtained if trained on clean preference data. We provide the detailed proof in Appendix A.

**Theorem 4.3** (Consistency of $\boldsymbol{p_\theta}$). *Given both the corrupted preference data distribution $\tilde{\mathcal{D}}$ induced by the flipping process and the unobserved clean preference data distribution $\mathcal{D}$, the following equality holds:*

$$\arg\min_\phi -\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \tilde{\mathcal{D}}}[\log \tilde{p}_\phi] = \arg\min_\phi -\mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}}[\log p_\phi],$$

*where $\tilde{p}_\phi = (1-\varepsilon_{\tilde{\boldsymbol{x}}})p_\phi + \varepsilon_{\tilde{\boldsymbol{x}}}(1-p_\phi)$ represents the predicted probability under the flipping model.*

This result holds under the specific parameterization where the flipping noise is marginalized via $\tilde{p}_\phi$. Furthermore, the core consistency principle holds equivalently under the DPO parameterization: replacing $p_\phi$ with $p_\theta$ and $\tilde{p}_\phi$

with $\tilde{p}_\theta$, the same result applies to the policy $\theta$ optimized via the DPO loss adjusted for flipping.

On the other hand, through a coordinate descent perspective on the iterative updates, the optimization of the flipping model reduces to the convex logistics regression. We present the following convergence result for self-containedness.

We first make the technical assumption that the features $h(\tilde{x})$ and the parameters $\omega$ for the preference flipping model are bounded, i.e., $\|h(\tilde{x})\| \le B$, $\|\omega\| \le B_\omega$. Then, we further assume that the constructed features satisfy the following coverage assumption, which is a common assumption in robust machine learning and the implied feature diversity is critical to prevent collapse or degenerate solutions.

**Assumption 4.4** (Feature coverage). Given the corrupted data distribution $\tilde{\mathcal{D}}$ and the feature map $h : \mathcal{X} \to \mathbb{R}^d$, the population covariance matrix of features satisfies
$$\lambda_{\min}\left(\mathbb{E}_{\tilde{x}\sim\tilde{\mathcal{D}}}\left[h(\tilde{x})h(\tilde{x})^\top\right]\right) > 0.$$

This assumption implies that the feature vectors are not concentrated in a low-dimensional subspace, which eliminates potential collinearity. The parameters $\omega$ are thus identifiable from the data.

Under the boundedness assumption, Assumption 4.4, and the condition that the reward or policy model provides accurate predictions ($p_\phi = p^*$ or $p_\theta = p^*$), we can establish fast convergence for the estimator of the flipping parameters $\omega$ obtained via gradient descent:

**Theorem 4.5** (Linear Convergence of $\hat{\omega}$). *Given loss function in Equation* (7), *if* $p_\phi = p^*$ *or* $p_\theta = p^*$ *and for the gradient descent update with step size* $\eta > 0$:
$$\omega^{(t+1)} = \omega^{(t)} - \eta\nabla_\omega\mathcal{L}_{FA\text{-}DPO}(\omega^{(t)}),$$
*the sequence of parameter estimates* $\{\omega^{(t)}\}$ *converges* Q-linearly *to the optimal parameter* $\omega^*$:
$$\|\omega^{(t+1)} - \omega^*\|^2 \le (1 - \eta\mu)\|\omega^{(t)} - \omega^*\|^2,$$
*for some* $\mu > 0$, *the convergence holds when* $0 < \eta < \frac{2}{L}$, *where $L$ is the smoothness constant for* $\mathcal{L}_{FA\text{-}DPO}$.

Theorem 4.5 provides the important guarantee that, when initialized with an accurate reward model or policy (reflecting the true clean preference probability $p^*$), the gradient descent update on the preference flipping model parameters $\omega$ converges rapidly at a linear rate. Consequently, this justifies the iterative procedures used in our algorithm.

# 5 Experiments

We investigate the effectiveness of FA-DPO mainly on aligning the LLMs with preference pairs under the instance-dependent preference flipping setting.

## 5.1 Experiment Setup

**Datasets and models** We conduct experiments mainly on two preference datasets: UltraFeedback (Cui et al. 2024) and Anthropic's HH_Golden (Anthropic 2024). For UltraFeedback, we follow the spirit of Chen et al. (2025) and first filter out samples with low score margins ($\le 0.5$) to obtain a cleaner subset for noise injection. This results in 42.4k/61.1k training samples and 1.4k/2k test samples. However, in experiments without manual label flipping, we retain the full training set while using the filtered test set for evaluation. For the backbone LLM models, we first test Pythia-1B on various levels of noise, and then we scale our method on larger models, LLama-3.1-8B and Mistral-7B, to validate the model's generation capabilities under data contamination. The detailed training setup are listed in Appendix E.

**Instance-dependent flipping** To validate our designs, we construct datasets with simulated instance-dependent preference flipping based on the following steps:

1. Randomly initialize a preference flipping model $\mathcal{N}_\varepsilon(\vartheta)$ (logistic model) on data features discussed in Section 4.3, and the model outputs a flip probability $\varepsilon_{\tilde{x}}$ for the current preference pair.

2. To determine whether to flip a preference sample, we set a threshold $\tau$ such that if the flipping probability $\varepsilon_{\tilde{x}}$ computed by the noise model is high enough to exceed this threshold, we would flip the preference. This is to maintain controlled stochasticity during preference flipping. During our experiment, we choose $\tau = 0.8$.

3. To control the total ratio of flips, we train the parameters $\vartheta$ on the training set $\mathcal{D}_{\text{train}}$ using the following loss function during noise model initialization, which forces the quantile of the flipping probability in the interval $[\tau, 1]$ to be around a pre-determined flipping ratio $\eta$.
$$\mathcal{L}_\varepsilon = \mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}_{\text{train}}}\left[\eta - \frac{\mathbb{I}\left(\mathcal{N}_\varepsilon(x,y_w,y_l;\vartheta) \ge \tau\right)}{|\mathcal{D}_{\text{train}}|}\right]^2.$$

After training on the contaminated training set, we test the model's performance on the clean test set which has no noise injected. For experiments with simulated flipping ($\eta > 0$), we intentionally restrict the capability of preference flipping model via only using length-based features, recovering the practical mismatch between estimations and real-world mechanisms. For clean datasets without added flipping ($\eta = 0$), we utilize the full feature set to model the original dataset's inherent noise patterns.

**Evaluation metrics** We use *prediction accuracy* (ACC) and *win rate* (WR) as our evaluation metrics.

In particular, we evaluate each model's *prediction accuracy* on a clean test set $\mathcal{D}_{\text{test}}$ across different noise levels. The accuracy is computed by comparing the predicted rewards for "chosen" versus "rejected" responses:
$$\text{Acc} = \mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}_{\text{test}}}\left[\mathbb{I}\left(\hat{r}_\theta(x,y_w) > \hat{r}_\theta(x,y_l)\right)\right],$$
where $\mathbb{I}$ is the indicator function, and $\hat{r}_\theta$ is the implicit reward model induced by DPO.

To compute the *win rate*, we employ different evaluators based on model capabilities: DeepSeek-V3 for Pythia-1B, and GPT-4o for both LLama-3.1-8B and Mistral-7B. The calculation is performed as follows:
$$\text{WR} = \frac{\#(\text{Win}) + \#(\text{Tie})/2}{\#(\text{Comparisons})},$$

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Ultrafeedback** (61.1k) | | | | | | | | | | |
| $\eta$ | 0% | | 10% | | 20% | | 30% | | 40% | |
| Methods | Acc ↑ | WR ↑ | Acc ↑ | WR ↑ | Acc ↑ | WR ↑ | Acc ↑ | WR ↑ | Acc ↑ | WR ↑ |
| DPO | 68.22 | **67.00** | 61.77 | 60.10 | 58.63 | 56.35 | 55.43 | 64.70 | 51.87 | 64.35 |
| SIMPO | 63.64 | 62.60 | 63.73 | 51.20 | 57.81 | 56.50 | 55.53 | 54.90 | 51.77 | 55.80 |
| ROPO | <u>70.75</u> | 66.80 | 64.93 | 66.85 | 58.93 | 65.55 | 55.23 | 64.40 | <u>54.07</u> | <u>67.25</u> |
| cDPO | 67.20 | 66.65 | 62.57 | <u>67.80</u> | 59.00 | <u>66.05</u> | 56.67 | <u>66.55</u> | 53.93 | 67.05 |
| rDPO | 70.13 | 65.60 | <u>65.90</u> | 56.65 | <u>61.70</u> | 54.85 | <u>56.87</u> | 57.85 | 47.67 | 57.80 |
| FA-DPO | **73.05** | <u>66.85</u> | **67.20** | **68.45** | **69.77** | **66.90** | **70.97** | **69.80** | **70.77** | **69.80** |
| (#Improv.) | +2.30 | -0.15 | +1.30 | +0.65 | +8.07 | +0.85 | +14.10 | +3.25 | +16.70 | +2.55 |
| **HH_Golden** (42.5k) | | | | | | | | | | |
| DPO | 98.89 | 87.00 | 93.50 | 34.05 | 83.53 | 29.95 | 73.30 | 21.05 | 58.63 | 26.85 |
| SIMPO | 98.27 | 65.10 | 93.63 | 29.95 | 82.83 | 24.90 | 73.33 | 22.40 | 60.93 | 27.50 |
| ROPO | 98.80 | 70.50 | 93.27 | 81.15 | 82.87 | 54.30 | 72.33 | 53.50 | <u>62.00</u> | 49.20 |
| cDPO | <u>99.24</u> | <u>88.35</u> | 93.20 | 56.95 | 82.50 | <u>62.30</u> | 73.30 | 37.75 | 60.87 | 49.05 |
| rDPO | 96.80 | 77.30 | <u>96.67</u> | <u>82.80</u> | <u>87.83</u> | 53.55 | <u>73.73</u> | <u>58.10</u> | 56.80 | <u>54.10</u> |
| FA-DPO | **99.61** | **88.70** | **99.17** | **88.90** | **98.10** | **82.85** | **99.02** | **87.70** | **98.83** | **78.60** |
| (#Improv.) | +0.37 | +0.35 | +2.50 | +6.10 | +10.27 | +20.55 | +25.29 | +29.60 | +36.83 | +24.50 |

Table 1: Model performance under different preference flipping ratios (0%-40%) for Pythia-1B on Ultrafeedback and HH_Golden. All values represent percentages, with **bold** indicating the highest score per column and <u>underline</u> denoting the runner-up. The (#Improv.) row quantifies absolute performance gain over the runner-up baseline.

where #(Win) and #(Tie) represent the number of wins and ties compared to the *reference model* (In our case, the SFT model), respectively, and #(Comparisons) denotes the total number of comparisons between the two models.

For *prediction accuracy*, we compute 5 runs for each method and report the mean value, with each run extracting randomly 1k samples from $\mathcal{D}_{\text{test}}$ for testing. For *win rate*, we extract 1k prompts from $\mathcal{D}_{\text{test}}$, and generate 1k responses from the trained policy model and the corresponding SFT model, respectively. We refer the readers to Appendix E for the detailed prompt template of LLM evaluators.

## 5.2 Results

**Discriminative performance** The prediction accuracy, derived from the policy model's implicit reward signals, reflects the model's discriminative performance in distinguishing between chosen and rejected responses. From Table 1 and Table 2, we can observe that, as the total flipping ratio increases, all baselines experience sharp prediction accuracy decreases across both datasets. Among the two datasets, HH_Golden exhibits a sharper decline than Ultrafeedback due to its larger chosen-rejected gap, also explaining why all methods achieve much higher accuracy on clean HH_Golden versus Ultrafeedback. Among all baselines, DPO and SIMPO show the highest sensitivity to preference flipping, while the robust methods (ROPO, cDPO, rDPO) exhibit moderate resilience. FA-DPO demonstrates superior robustness across all flipping ratios, datasets and models, consistently outperforming other approaches.

**Generative performance** The generation performance shows a similar trend with the discriminative performance. As shown in Table 1, DPO and SIMPO exhibit rapid performance drops with increasing flipping ratios, revealing their inherent fragility to preference noise. This vulnera-

bility is pronounced when chosen and rejected responses come from distinct distributions (like in HH_Golden), reinforcing the overfitting patterns in these methods. FA-DPO proves effective for preserving the generation capabilities of LLMs under instance-dependent preference flipping scenarios, with consistent performance across all conditions, as well as maintaining performance under the clean dataset.

**Preference flipping model** Then we compare the flipping probabilities predicted by the learned model with the ground-truth flipping distribution. Predicted flipping probabilities show close alignment with actual flipping distribution, as shown in Figure 1 (a). Moreover, Figure 1 (b) reveals a clear separation of the predicted flipping probabilities between flipped and non-flipped samples, indicating distinct decision boundaries between the two kinds of samples. We also plot the relationship between the flipping probabilities with length-based feature in Figure 1 (c) to show that our model captures the relation between noise and data features.

## 5.3 Ablation Studies

**Hyperparameters in iterative training** We conduct an ablation study on the hyperparameters in iterative training, first examining the impact of the *warmup* operation, then evaluating different iteration step combinations under a 20% flipping ratio on HH_Golden. Table 3 shows that *warmup* plays a vital role in performance improvement. We remark that the preference flipping model needs to be sufficiently trained to guide the policy learning, especially when the initial policy model is not well trained. When the flipping model performs well, training more policy steps yields improved performance.

**Computational cost** Although FA-DPO requires training an auxiliary flipping model, we note the overall com-

(a) Prediction correlation  (b) Noise distribution by flipping status  (c) Length-based noise patterns
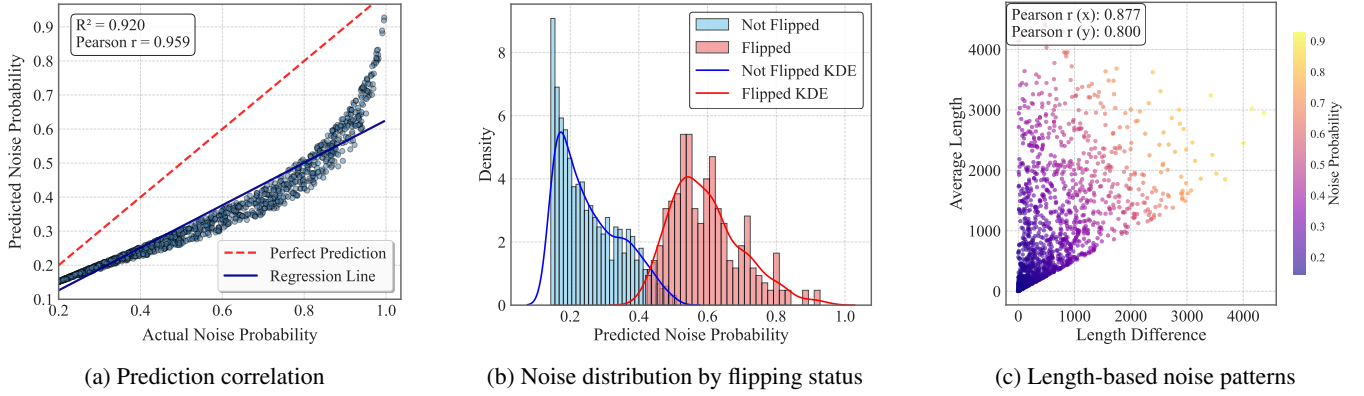
Figure 1: Characterization of learned preference flipping distribution. (a) Correlation between actual and predicted noise probabilities with regression line; (b) Predicted flipping distributions separated by flipping status; (c) Pattern of predicted flipping distribution with length-based features.

| | **LLama-3.1-8B** | | | |
|---|---|---|---|---|
| $\eta$ | 20% | | 40% | |
| Methods | Acc ↑ | WR ↑ | Acc ↑ | WR ↑ |
| DPO | 73.89 | 62.90 | 64.96 | <u>56.10</u> |
| SIMPO | 66.59 | 62.80 | 57.07 | 38.95 |
| ROPO | <u>75.82</u> | 62.75 | <u>66.59</u> | 50.25 |
| cDPO | 75.22 | <u>64.50</u> | 65.70 | <u>56.10</u> |
| rDPO | 73.96 | 59.75 | 63.99 | 55.90 |
| FA-DPO | **78.80** | **65.10** | **78.87** | **68.50** |
| (#Improv.) | +2.98 | +0.60 | +12.28 | +12.40 |
| | **Mistral-7B** | | | |
| DPO | 71.35 | 45.75 | 62.05 | 35.35 |
| SIMPO | 67.04 | <u>54.75</u> | 63.77 | <u>45.00</u> |
| ROPO | <u>73.96</u> | 54.30 | <u>64.36</u> | 36.35 |
| cDPO | 72.99 | 47.75 | 63.24 | 34.60 |
| rDPO | 71.95 | 44.75 | 63.10 | 36.65 |
| FA-DPO | **78.05** | **61.35** | **78.49** | **59.00** |
| (#Improv.) | +4.09 | +6.60 | +14.13 | +14.00 |

Table 2: Model performance under flipping ratios (20% and 40%) for LLama-3.1-8B and Mistral-7B on Ultrafeedback.

| Warmup | Iteration Batches | | Metrics | |
|---|---|---|---|---|
| | Noise | Policy | Acc↑ | WR↑ |
| | 20 | 20 | 84.96 | 57.35 |
| No | 20 | 50 | 77.00 | 53.20 |
| | 50 | 50 | 83.16 | <u>76.85</u> |
| | 20 | 20 | **98.56** | 70.04 |
| Yes | 20 | 50 | <u>98.40</u> | **82.90** |
| | 50 | 50 | 96.80 | **82.90** |

Table 3: Ablation studies on warmup and policy/noise iteration steps.

putational cost remains comparable to—or even lower than—that of standard DPO. This is because we limit training to a single epoch, keeping the total data usage consistent across all methods. Additionally, the PPL and reward margin features required for training the preference flipping model are derived directly from the policy's forward pass log-likelihoods, introducing no extra computation overhead.

## 6 Conclusion

In this paper, we tackle the challenging scenario of instance-dependent noisy human feedback through introducing a framework that simultaneously models preference flipping and post-trains the LLM, on the basis of the RLHF and the DPO frameworks. Instead of directly modeling the noise inside the BT model, our approach separates human intention from the noising process by assuming a post-transition after the BT model forming the preference with encoded probabilities. This process features the stochastic transformation from groundtruth to noisy labels, providing a more realistic representation. By integrating the MLE-based BT model with preference flipping probabilities, we can then learn a statistically consistent estimator. In more detail, our algorithm iteratively updates the noise model and fine-tunes the LLM parameters, and our implementation on DPO is achieved with little additional resource consumption. Empirically, our approach adopts relevant sequence features to model flipping ratios and yields high probability for preference flipping, as expected. In evaluations on instance-dependent noisy human preference datasets, our algorithm demonstrates higher predictive accuracy compared to vanilla DPO and other baseline methods.

# References

Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Anthropic. 2024. Anthropic's HH_Golden.

Azar, M. G.; Guo, Z. D.; Piot, B.; Munos, R.; Rowland, M.; Valko, M.; and Calandriello, D. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, 4447–4455. PMLR.

Bagnell, J. A. 2005. Robust supervised learning. In *AAAI*, 714–719.

Bai, Y.; Jones, A.; Ndousse, K.; Askell, A.; Chen, A.; Das-Sarma, N.; Drain, D.; Fort, S.; Ganguli, D.; Henighan, T.; et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Bradley, R. A.; and Terry, M. E. 1952. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika*, 39(3/4): 324–345.

Bukharin, A.; Hong, I.; Jiang, H.; Li, Z.; Zhang, Q.; Zhang, Z.; and Zhao, T. 2024. Robust reinforcement learning from corrupted human feedback. *Advances in Neural Information Processing Systems*, 37: 124093–124113.

Casper, S.; Davies, X.; Shi, C.; Gilbert, T. K.; Scheurer, J.; Rando, J.; Freedman, R.; Korbak, T.; Lindner, D.; Freire, P.; et al. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*.

Chakraborty, S.; Qiu, J.; Yuan, H.; Koppel, A.; Huang, F.; Manocha, D.; Bedi, A. S.; and Wang, M. 2024. MaxMin-RLHF: Towards equitable alignment of large language models with diverse human preferences. *arXiv preprint arXiv:2402.08925*.

Chen, G. H.; Chen, S.; Liu, Z.; Jiang, F.; and Wang, B. 2024. Humans or llms as the judge? a study on judgement biases. *arXiv preprint arXiv:2402.10669*.

Chen, P.; Chen, X.; Yin, W.; and Lin, T. 2025. ComPO: Preference alignment via comparison oracles. *arXiv preprint arXiv:2505.05465*.

Cheng, J.; Xiong, G.; Dai, X.; Miao, Q.; Lv, Y.; and Wang, F.-Y. 2024. RIME: Robust Preference-based Reinforcement Learning with Noisy Preferences. *arXiv preprint arXiv:2402.17257*.

Chowdhury, S. R.; Kini, A.; and Natarajan, N. 2024. Provably robust dpo: Aligning language models with noisy feedback. *arXiv preprint arXiv:2403.00409*.

Coste, T.; Anwar, U.; Kirk, R.; and Krueger, D. 2023. Reward model ensembles help mitigate overoptimization. *arXiv preprint arXiv:2310.02743*.

Cui, G.; Yuan, L.; Ding, N.; Yao, G.; He, B.; Zhu, W.; Ni, Y.; Xie, G.; Xie, R.; Lin, Y.; et al. 2024. ULTRAFEEDBACK: Boosting Language Models with Scaled AI Feedback. In *Forty-first International Conference on Machine Learning*.

Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Yang, A.; Fan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Ethayarajh, K.; Xu, W.; Muennighoff, N.; Jurafsky, D.; and Kiela, D. 2024. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*.

Fisch, A.; Eisenstein, J.; Zayats, V.; Agarwal, A.; Beirami, A.; Nagpal, C.; Shaw, P.; and Berant, J. 2024. Robust preference optimization through reward model distillation. *arXiv preprint arXiv:2405.19316*.

Gao, Y.; Alon, D.; and Metzler, D. 2024. Impact of preference noise on the alignment performance of generative language models. *arXiv preprint arXiv:2404.09824*.

Go, D.; Korbak, T.; Kruszewski, G.; Rozen, J.; Ryu, N.; and Dymetman, M. 2023. Aligning language models with preferences through f-divergence minimization. *arXiv preprint arXiv:2302.08215*.

Hendrycks, D.; Mazeika, M.; Kadavath, S.; and Song, D. 2019. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32.

Hong, J.; Lee, N.; and Thorne, J. 2024. Orpo: Monolithic preference optimization without reference model. *arXiv preprint arXiv:2403.07691*, 2(4): 5.

Huang, K.; Wu, J.; Chen, Z.; Wang, X.; Gao, J.; Ding, B.; Wu, J.; He, X.; and Wang, X. 2025. Larger or Smaller Reward Margins to Select Preferences for LLM Alignment? In *Forty-second International Conference on Machine Learning*.

Kim, T.; Ko, J.; Choi, J.; Yun, S.-Y.; et al. 2021. Fine samples for learning with noisy labels. *Advances in Neural Information Processing Systems*, 34: 24137–24149.

Kong, K.; Xu, X.; Wang, D.; Zhang, J.; and Kankanhalli, M. S. 2024. Perplexity-aware correction for robust alignment with noisy preferences. *Advances in Neural Information Processing Systems*, 37: 28296–28321.

Lee, K.; Smith, L.; Dragan, A.; and Abbeel, P. 2021. B-pref: Benchmarking preference-based reinforcement learning. *arXiv preprint arXiv:2111.03026*.

Liang, X.; Chen, C.; Wang, J.; Wu, Y.; Fu, Z.; Shi, Z.; Wu, F.; and Ye, J. 2024. Robust preference optimization with provable noise tolerance for llms. *arXiv preprint arXiv:2404.04102*.

Liu, Y.; Cheng, H.; and Zhang, K. 2023. Identifiability of label noise transition matrix. In *International Conference on Machine Learning*, 21475–21496. PMLR.

Meng, Y.; Xia, M.; and Chen, D. 2024. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.

Mitchell, E. 2023. A note on DPO with noisy preferences & relationship to IPO.

Muslea, I.; Minton, S.; and Knoblock, C. A. 2002. Active+ semi-supervised learning= robust multi-view learning. In *ICML*, volume 2, 435–442. Citeseer.

Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.

Patrini, G.; Rozza, A.; Krishna Menon, A.; Nock, R.; and Qu, L. 2017. Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1944–1952.

Peng, X. B.; Kumar, A.; Zhang, G.; and Levine, S. 2019. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*.

Rafailov, R.; Sharma, A.; Mitchell, E.; Manning, C. D.; Ermon, S.; and Finn, C. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Ramesh, S. S.; Hu, Y.; Chaimalas, I.; Mehta, V.; Sessa, P. G.; Ammar, H. B.; and Bogunovic, I. 2024. Group Robust Preference Optimization in Reward-free RLHF. *arXiv preprint arXiv:2405.20304*.

Reid, M. D.; and Williamson, R. C. 2010. Composite binary losses. *The Journal of Machine Learning Research*, 11: 2387–2422.

Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Song, H.; Kim, M.; Park, D.; Shin, Y.; and Lee, J.-G. 2022. Learning from noisy labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning systems*, 34(11): 8135–8153.

Wu, J.; Xie, Y.; Yang, Z.; Wu, J.; Chen, J.; Gao, J.; Ding, B.; Wang, X.; and He, X. 2024a. Towards Robust Alignment of Language Models: Distributionally Robustifying Direct Preference Optimization. *arXiv preprint arXiv:2407.07880*.

Wu, J.; Xie, Y.; Yang, Z.; Wu, J.; Gao, J.; Ding, B.; Wang, X.; and He, X. 2024b. *beta*-DPO: Direct Preference Optimization with Dynamic *beta*. *Advances in Neural Information Processing Systems*, 37: 129944–129966.

Wu, Z.; Hu, Y.; Shi, W.; Dziri, N.; Suhr, A.; Ammanabrolu, P.; Smith, N. A.; Ostendorf, M.; and Hajishirzi, H. 2024c. Fine-grained human feedback gives better rewards for language model training. *Advances in Neural Information Processing Systems*, 36.

Xia, X.; Liu, T.; Han, B.; Gong, M.; Yu, J.; Niu, G.; and Sugiyama, M. 2021. Sample selection with uncertainty of losses for learning with noisy labels. *arXiv preprint arXiv:2106.00445*.

Xia, X.; Liu, T.; Han, B.; Wang, N.; Gong, M.; Liu, H.; Niu, G.; Tao, D.; and Sugiyama, M. 2020. Part-dependent label noise: Towards instance-dependent label noise. *Advances in Neural Information Processing Systems*, 33: 7597–7610.

Xiao, J.; Li, Z.; Xie, X.; Getzen, E.; Fang, C.; Long, Q.; and Su, W. J. 2024. On the Algorithmic Bias of Aligning Large Language Models with RLHF: Preference Collapse and Matching Regularization. *arXiv preprint arXiv:2405.16455*.

Yuan, Z.; Yuan, H.; Tan, C.; Wang, W.; Huang, S.; and Huang, F. 2023. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*.

Zhang, C.; Bengio, S.; Hardt, M.; Recht, B.; and Vinyals, O. 2021. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3): 107–115.

Zhang, C.; Vinyals, O.; Munos, R.; and Bengio, S. 2018. A study on overfitting in deep reinforcement learning. *arXiv preprint arXiv:1804.06893*.

Zhang, S.; Chen, Z.; Chen, S.; Shen, Y.; Sun, Z.; and Gan, C. 2024. Improving reinforcement learning from human feedback with efficient reward model ensemble. *arXiv preprint arXiv:2401.16635*.

Zhao, Y.; Joshi, R.; Liu, T.; Khalman, M.; Saleh, M.; and Liu, P. J. 2023. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.

Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing, E.; et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36: 46595–46623.

Ziegler, D. M.; Stiennon, N.; Wu, J.; Brown, T. B.; Radford, A.; Amodei, D.; Christiano, P.; and Irving, G. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.

# Technical Appendix

## A    Theoretical Proofs

### A.1    Proof for Theorem 4.2

*Proof.* We derive the gradient analysis for our proposed FA-DPO objective, highlighting its connection to standard DPO. The loss function for FA-DPO is defined as:

$$\mathcal{L}_{\text{FA-DPO}} = -\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \tilde{\mathcal{D}}} \left[ \log \left( (1 - \varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta + \varepsilon_{\tilde{\boldsymbol{x}}} (1 - p_\theta) \right) \right],$$

where $p_\theta$ represents the policy's preference probability:

$$p_\theta = \sigma \left( \beta \log \frac{\pi_\theta(\tilde{y}_w \mid x)}{\pi_{\text{ref}}(\tilde{y}_w \mid x)} - \beta \log \frac{\pi_\theta(\tilde{y}_l \mid x)}{\pi_{\text{ref}}(\tilde{y}_l \mid x)} \right).$$

This formulation explicitly incorporates the *instance-dependent* flipping rate $\varepsilon_{\tilde{\boldsymbol{x}}}$, which distinguishes FA-DPO from methods assuming fixed noise rates.

We first compute the gradient of $p_\theta$ with respect to the policy parameters $\theta$. Applying the chain rule and sigmoid derivative $\nabla \sigma(z) = \sigma(z)(1 - \sigma(z))$ yields:

$$\nabla_\theta p_\theta = -\beta \cdot p_\theta (1 - p_\theta) \left( \nabla_\theta \log \pi_\theta(\tilde{y}_w \mid x) - \nabla_\theta \log \pi_\theta(\tilde{y}_l \mid x) \right).$$

This expression shares the same structure as DPO's gradient but will be scaled differently in our loss.

Now, taking the gradient of the full FA-DPO objective:

$$\nabla_\theta \mathcal{L}_{\text{FA-DPO}} = -\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \tilde{\mathcal{D}}} \left[ \frac{1}{(1 - \varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta + \varepsilon_{\tilde{\boldsymbol{x}}} (1 - p_\theta)} \cdot \nabla_\theta \left( (1 - \varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta + \varepsilon_{\tilde{\boldsymbol{x}}} (1 - p_\theta) \right) \right]$$

$$= -\beta \cdot \mathbb{E} \left[ \frac{(1 - 2\varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta}{(1 - 2\varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta + \varepsilon_{\tilde{\boldsymbol{x}}}} \cdot (1 - p_\theta) \left( \nabla_\theta \log \pi_\theta(\tilde{y}_w \mid x) - \nabla_\theta \log \pi_\theta(\tilde{y}_l \mid x) \right) \right].$$

From this derivation, we identify the weighting coefficient for FA-DPO as:

$$\zeta_{\text{FA-DPO}} = \frac{(1 - 2\varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta}{(1 - 2\varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta + \varepsilon_{\tilde{\boldsymbol{x}}}} \cdot (1 - p_\theta).$$

To understand how FA-DPO relates to existing methods, recall the weighting coefficients derived by Chowdhury, Kini, and Natarajan (2024) for constant noise rates:

$$\zeta_{\text{DPO}} = 1 - p_\theta$$
$$\zeta_{\text{cDPO}} = \zeta_{\text{DPO}} - \varepsilon$$
$$\zeta_{\text{rDPO}} = \zeta_{\text{DPO}} + \frac{\varepsilon}{1 - 2\varepsilon},$$

where $\varepsilon$ denotes a *fixed* flipping rate. Comparing these expressions, we observe that FA-DPO's weight adaptively adjusts $\zeta_{\text{DPO}}$ by:

$$\zeta_{\text{FA-DPO}} = \underbrace{\frac{(1 - 2\varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta}{(1 - 2\varepsilon_{\tilde{\boldsymbol{x}}}) p_\theta + \varepsilon_{\tilde{\boldsymbol{x}}}}}_{\text{instance-dependent scaling}} \cdot \zeta_{\text{DPO}},$$

which completes the proof.

$\square$

## A.2 Proof for Theorem 4.3

*Proof.* For any corrupted preference sample $\tilde{x} = (x, \tilde{y}_w, \tilde{y}_l)$, we express the FA-DPO loss in vector form:

$$\ell_{\text{FA-DPO}}(\tilde{x}) = -\log\left(\Lambda(\tilde{x})p_\theta\right),$$

where $p_\theta = [p_\theta, 1 - p_\theta]^\top$ is the preference probability vector, and the flipping matrix $\Lambda(\tilde{x})$ is defined as:

$$\Lambda(\tilde{x}) = \begin{bmatrix} 1 - \varepsilon_{\tilde{x}} & \varepsilon_{\tilde{x}} \\ \varepsilon_{\tilde{x}} & 1 - \varepsilon_{\tilde{x}} \end{bmatrix}.$$

For all $\varepsilon_{\tilde{x}} \neq 0.5$, $\Lambda(\tilde{x})$ is invertible, with determinant $\det(\Lambda) = (1 - \varepsilon_{\tilde{x}})^2 - \varepsilon_{\tilde{x}}^2 = 1 - 2\varepsilon_{\tilde{x}}$.

Define the logit vector $h(\tilde{x})$ as:

$$h(\tilde{x}) = \left[\beta\log\frac{\pi_\theta(\tilde{y}_w \mid x)}{\pi_{\text{ref}}(\tilde{y}_w \mid x)} - \beta\log\frac{\pi_\theta(\tilde{y}_l \mid x)}{\pi_{\text{ref}}(\tilde{y}_l \mid x)}, \beta\log\frac{\pi_\theta(\tilde{y}_l \mid x)}{\pi_{\text{ref}}(\tilde{y}_l \mid x)} - \beta\log\frac{\pi_\theta(\tilde{y}_w \mid x)}{\pi_{\text{ref}}(\tilde{y}_w \mid x)}\right]^\top.$$

The preference probability is related to this logit vector through the sigmoid function: $p_\theta = \sigma(h(\tilde{x}))$.

By Reid and Williamson (Reid and Williamson 2010), the binary cross-entropy (BCE) loss function is proper composite, meaning it satisfies:

$$\ell(y, h) = \phi(\langle y, h\rangle) + c(h),$$

where $y$ is the one-hot encoded true label, and $\phi$ is a strictly convex function. This property holds for our FA-DPO loss formulation.

Given that (1) $\Lambda(\tilde{x})$ is invertible for $\varepsilon_{\tilde{x}} \neq 0.5$, and (2) FA-DPO retains the composite loss property through the linear transformation $\Lambda(\tilde{x})$, we directly apply Theorem 2 from Patrini et al. (Patrini et al. 2017):

$$\arg\min_{h} -\mathbb{E}_{\tilde{x}\sim\tilde{\mathcal{D}}}\left[\ell_{\text{FA-DPO}}\right] = \arg\min_{h} -\mathbb{E}_{x\sim\mathcal{D}}\left[\ell_{\text{DPO}}\right]$$

$$= \arg\min_{h} -\mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}}\left[\log\sigma\left(\beta\log\frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta\log\frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)}\right)\right],$$

which yields the desired result

$$\arg\min_{\phi} -\mathbb{E}_{\tilde{x}\sim\tilde{\mathcal{D}}}[\log\tilde{p}_\phi] = \arg\min_{\phi} -\mathbb{E}_{x\sim\mathcal{D}}[\log p_\phi].$$

$\square$

## A.3 Proof for Theorem 4.5

*Proof.* To establish convergence of the gradient descent updates to the true parameters $w^*$ of FA-DPO loss function, we first show that the gradient vanishes at optimum.

For a sample triplet $\tilde{x} = (x, \tilde{y}_w, \tilde{y}_l)$, we denote a random variable $\tilde{Y} = 1$ to represent $\tilde{y}_w \succ \tilde{y}_l \mid x$, similarly, $\tilde{Y} = 0$ to represent $\tilde{y}_w \prec \tilde{y}_l \mid x$. Then We write the population loss for FA-DPO:

$$\mathcal{L}_{\text{FA-DPO}} = -\mathbb{E}_{\tilde{x},\tilde{Y}\sim\mathcal{D}}\left[\mathbb{I}\{\tilde{Y} = 1\}\log\tilde{p} + \mathbb{I}\{\tilde{Y} = 0\}\log(1 - \tilde{p})\right],$$

where $\tilde{p} = (1 - \varepsilon(\tilde{x}; \omega))p + \varepsilon(\tilde{x}; \omega)p$.

We consider the parameterization of $p$ as $p_\phi$ in the standard reward learning paradigm, and the derived results apply equivalently to that of the DPO parameterization ($p_\theta$).

We derive the gradient of $\ell_{\text{FA-DPO}}$ for sample $\tilde{x}$ with respect to $\omega$ as:

$$\nabla_\omega \ell(\tilde{x}; \omega) = -(1 - 2p_\phi)\left[\frac{1}{\tilde{p}_\phi} \cdot \mathbb{I}\{\tilde{Y} = 1\} + \frac{1}{1 - \tilde{p}_\phi} \cdot \mathbb{I}\{\tilde{Y} = 0\}\right]\nabla_\omega \varepsilon(\tilde{x}; \omega).$$

Based on the assumption that $p_\phi = p^*$, at the point $\omega^*$, we have

$$\mathbb{E}_{\tilde{x},\tilde{Y}\sim\mathcal{D}}[\mathbb{I}\{\tilde{Y} = 1\} \mid \tilde{x}] = \tilde{p}^*.$$

Consequently, the expected conditional partial derivative vanishes:

$$\mathbb{E}\left[\nabla_\omega \ell_{\text{FA-DPO}}(\tilde{\boldsymbol{x}}; \omega) \mid \tilde{\boldsymbol{x}}\right] = 0$$

According to the law of total expectation, we get the following equation:

$$\mathbb{E}[\mathbb{E}[\nabla_\omega \ell_{\text{FA-DPO}} \mid \tilde{\boldsymbol{x}}]] = \mathbb{E}[\nabla_\omega \ell_{\text{FA-DPO}}] = 0$$

Then we further compute the Hessian of $\mathcal{L}_{\text{FA-DPO}}$ with respect to $\omega$. Applying the same trick in gradient derivation, the second-order terms of $\nabla_\omega \varepsilon(\omega)$ in Hessian are canceled out on the conditional expectation, therefore, we have the following result:

$$\nabla_\omega^2 \mathcal{L}_{\text{FA-DPO}} = \mathbb{E}\left[(1-2p)^2 \left(\frac{1}{\tilde{p}^*} + \frac{1}{1-\tilde{p}^*}\right) (\nabla_\omega \varepsilon(\tilde{\boldsymbol{x}}; \omega))(\nabla_\omega \varepsilon(\tilde{\boldsymbol{x}}; \omega))^\top\right]$$

Combined with $\sigma'(\langle \omega^*, z_i \rangle) \geq c_\sigma > 0$ and Theorem 4.4, we have that

$$\mathbb{E}[\nabla_\omega^2 \mathcal{L}_{\text{FA-DPO}}] \succeq \frac{4\delta^2 c_\sigma^2}{1-\delta} \cdot \mathbb{E}[h(\tilde{\boldsymbol{x}})h(\tilde{\boldsymbol{x}})^\top] \succeq \frac{4\delta^2 c_\sigma^2 \gamma}{1-\delta} I,$$

where $|1-2p| \geq \delta > 0$. Therefore, Hessian satisfies $\nabla_w^2 \mathcal{L} \succeq \mu I$ where $\mu = \frac{4\delta^2 c_\sigma^2 \gamma}{1-\delta}$.

Lipschitz smoothness follows from bounded parameters, as we have assumed that features satisfy $\|h(\tilde{\boldsymbol{x}})h(\tilde{\boldsymbol{x}})\| \leq B_z$, parameters $\|w\| \leq B_w$, as the sigmoid derivatives are bounded ($|\sigma'| \leq \frac{1}{4}$, $|\sigma''| \leq \frac{1}{4\sqrt{3}}$), Consequently, $\|\nabla_w^2 \mathcal{L}_{\text{FA-DPO}}\|_2 \leq L < \infty$ globally.

Finally, with $\mathcal{L}$ being $\mu$-strongly convex and $L$-smooth near $w^*$ and $\nabla_w \mathcal{L}(w^*) = \boldsymbol{0}$, gradient descent with step size $\eta_t < 2/L$ converges linearly:

$$\|w^{(t+1)} - w^*\|_2^2 \leq (1 - \eta_t \mu)\|w^{(t)} - w^*\|_2^2.$$

If the initial parameters are properly configured via bounded distance $\|w^{(0)} - w^*\|$. Setting $\rho = \sqrt{1-\eta\mu} < 1$ yields:

$$\|w^{(t)} - w^*\|_2 \leq \rho^t \|w^{(0)} - w^*\|_2.$$

$\square$

# B    Related Works

We review the previous works on preference alignment and (generalized) robust RLHF in this section, for the reader's convenience.

**Preference alignment.**    The most well-known approach for preference alignment is Reinforcement Learning from Human Feedback (Ziegler et al. 2019; Ouyang et al. 2022, RLHF), which involves training a *reward model* to capture human preferences and then guiding LLMs to *generate high-reward responses using reinforcement learning algorithms* such as Proximal Policy Optimization (PPO) (Schulman et al. 2017). However, in practice, RL-based methods can be complex and unstable during training (Rafailov et al. 2024; Wu et al. 2024c; Yuan et al. 2023). As a result, recent research has focused on simpler and more stable alternatives to RLHF (Rafailov et al. 2024; Zhao et al. 2023; Ethayarajh et al. 2024; Azar et al. 2024; Hong, Lee, and Thorne 2024; Meng, Xia, and Chen 2024).

Among these, a promising direction is to use contrastive or ranking loss to adjust the likelihood of output sequences. Specifically, RRHF (Yuan et al. 2023) introduces a ranking loss to increase the likelihood for better responses and decrease it for worse ones. Sequence Likelihood Calibration (SLiC) (Zhao et al. 2023) uses a range of calibration losses to align the model outputs with reference sequences in the latent space thereof. Additionally, Direct preference optimization (DPO) (Rafailov et al. 2024) offers an important approach by implicitly optimizing the same objective as existing RLHF methods, enabling human preference alignment directly through a simple cross-entropy loss. Due to the simplicity of DPO, a flurry of subsequent algorithms have introduced variants from different perspectives. For instance, SimPO (Meng, Xia, and Chen 2024) leverages the average log probability of a sequence as an implicit reward, removing the need for a reference model. And ORPO (Hong, Lee, and Thorne 2024) extends supervised fine-tuning (SFT) in preference alignment by employing an odds ratio to contrast favored and disfavored responses. Meanwhile, preference alignment methods without the reward model like $\Psi$PO (Azar et al. 2024) propose objectives based directly on pairwise preferences, bypassing the need for approximations typically used in constructing the reward model.

**RLHF against perturbations.** Most existing robust RLHF methods against perturbations are based on robust learning techniques from supervised learning (Bagnell 2005; Hendrycks et al. 2019; Muslea, Minton, and Knoblock 2002). Current approaches can be categorized into three main types. ❶ Noise fitting (Bukharin et al. 2024) involves making assumptions on the noise in the data and incorporating this modeling into the reward learning process, which will be jointly optimized with the parameterized reward function. The assumptions on the noise model significantly impact the performance of these algorithms. Various noise models have been proposed within the Bradley-Terry framework (Bradley and Terry 1952; Lee et al. 2021; Gao, Alon, and Metzler 2024); however, in the context of LLM alignment, only random flipping (Chowdhury, Kini, and Natarajan 2024) and sparse noise in the reward model (Bukharin et al. 2024) are typically considered, since the true reward model reflecting human intention is generally unknown.

In supervised learning, it has been observed that neural networks initially fit the clean data in the early stages of training and gradually overfit to noise (Zhang et al. 2021, 2018). Therefore, ❷ sample selection (or sample re-weighting) methods (Cheng et al. 2024) leverage this phenomenon to identify clean samples based on the loss values observed during the training process. Although empirically effective, these methods can mistakenly filter out true samples, thereby reducing the overall utility of the data (Xia et al. 2021; Kim et al. 2021). Some approaches address corruption in input data through ❸ robust loss design (Chowdhury, Kini, and Natarajan 2024; Liang et al. 2024), focusing on constructing loss functions that are more resistant to data noise; this technique is also known as label smoothing (Mitchell 2023).

**Generalized robust RLHF** Beyond robust approaches to mitigate data noise, generalized robust RLHF has incorporated a range of methods aimed at enhancing *resilience against various uncertainties*. Regularization techniques, for instance, have been employed to counteract overfitting and improve generalization within RLHF settings (Go et al. 2023; Xiao et al. 2024; Chowdhury, Kini, and Natarajan 2024). Other methods focus on bolstering reward model robustness; representative approaches, such as reward ensemble and distillation (Fisch et al. 2024; Coste et al. 2023; Zhang et al. 2024), help ensure more reliable feedback integration using multiple reward models to handle diverse data distributions. Other than ensemble, distributional robust optimization (DRO) (Wu et al. 2024a) has emerged as another robust framework, offering safeguards against distributional shifts in training data. Researchers also investigate diverse opinions among different groups of annotators; Ramesh et al. (2024); Chakraborty et al. (2024) proposed group robustness methodologies to address performance disparities across diverse subgroups, promoting fairness and equity. Together, these methods contribute to a more resilient and balanced RLHF framework, each addressing different facets of uncertainty while collectively enhancing robustness.

## C    Preliminaries

In this section, we introduce the basics of RLHF for LLM alignment.

### C.1   RLHF for LLM alignment

Let an LLM take an input (prompt) $x \in \mathcal{X}$ and generate an output (response) $y \in \mathcal{Y}$. In the contextual bandit formulation for RLHF, the LLM is viewed as a policy $\pi_\theta(y \mid x)$ parameterized by $\theta$, which outputs an action $y$ (response) based on the state $x$ (prompt).

The objective of LLM alignment is to optimize $\theta$ so that the output responses of the LLM align with human intentions. To represent human intentions, preference data with human annotations are collected for policy training. A preference data pair is collected in the form of $(y_1, y_2) \sim \pi_{\text{ref}}(y \mid x)$, where $\pi_{\text{ref}}(\cdot \mid \cdot)$ is a reference policy (detailed in the next paragraph). The preference data is further annotated by human labelers, denoted as $y_w \succ y_l \mid x$, where $y_w$ is the preferred response and $y_l$ is the dispreferred one in $(y_1, y_2)$ for the prompt $x$. Notably, the randomness in the preference dataset $D = \{(x^i, y_w^i, y_l^i)\}$ is two-fold: the responses $y_w^i, y_l^i$ are randomly generated, and the preference $\{y_w^i \succ y_l^i\}$ is a random event as well.

The standard RLHF pipeline consists of three stages (Ziegler et al. 2019). ❶ In the first stage, the pre-trained model undergoes one round of supervised fine-tuning (SFT) using a specific dataset for alignment, resulting in a so-called reference model $\pi_{\text{ref}}$. ❷ The second stage is reward modeling, where the Bradley-Terry (BT) model (Bradley and Terry 1952) is employed to connect the preference data $\{y_w^i \succ y_l^i\}$ to a reward model $r(x, y)$. The connection is formulated as:

$$p^*(y_w \succ y_l \mid x) = \sigma\big(r^*(x, y_w) - r^*(x, y_l)\big), \tag{11}$$

where $\sigma$ is the standard sigmoid function, and $r^*(\cdot)$ is the optimal reward model. Using the Bradley-Terry model described in Equation (11), the loss for learning the reward model is:

$$\begin{aligned} \mathcal{L}_R(\phi) &= -\mathbb{E}\left[\log p(y_w \succ y_l \mid x)\right] \\ &= -\mathbb{E}\left[\log \sigma\big(r_\phi(x, y_w) - r_\phi(x, y_l)\big)\right], \end{aligned} \tag{12}$$

where $r_\phi(\cdot)$ is the reward model parameterized by $\phi$, and the expectation is taken over $(x, y_w, y_l) \sim D$.

❸ After obtaining the reward model $r_\phi(\cdot)$, the third stage involves using reinforcement learning (RL) to optimize the LLM $\pi_\theta$ with the reward signals provided by $r_\phi(\cdot)$. The optimization objective for $\pi_\theta$ is formulated as:

$$
\begin{aligned}
\mathcal{L}_\pi(\theta) = &- \mathbb{E}_{x \sim \mathbb{P}_x, y \sim \pi_\theta} [r_\phi(x, y)] \\
&+ \beta \mathbb{D}_{KL} [\pi_\theta(y|x) || \pi_{\text{ref}}(y|x)],
\end{aligned}
\tag{13}
$$

where $\mathbb{P}_x$ represents the marginal distribution of the prompt $x$. The first term corresponds to reward maximization in standard RL optimization. The second term is a KL divergence that constrains the update of $\pi_\theta$ to not deviate from the reference model $\pi_{\text{ref}}$, mitigating the risk of out-of-distribution issues for the reward model and preventing mode collapse in the generation process; $\beta$ is the coefficient that weights the KL divergence between $\pi_\theta$ and $\pi_{\text{ref}}$.

## C.2 DPO as efficient RLHF

Direct preference learning (DPO) is an algorithm proposed for practical efficiency (compared to original RLHF), that merges the last two stages in RLHF, ❷ reward modeling and ❸ policy optimization, into a single step. In DPO, the policy is optimized directly using the offline dataset $D$ without constructing an explicit reward model.

As shown by Peng et al. (2019), the closed-form solution of the conditional distribution $\pi(y \mid x)$ that minimizes Equation (3) is:

$$
\pi_r(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r_\phi(x, y)\right),
\tag{14}
$$

where $Z(x) = \sum_y \pi_{\text{ref}}(y \mid x) \exp\left(1/\beta \cdot r_\phi(x, y)\right)$ is the partition function for normalization. Therefore, the information from reward model $r_\phi(x, y)$ can be implicitly recovered by the conditional density (equation 4). By applying this result to Equation (2), the consequent DPO loss function that includes only the parameterized $\pi_\theta$ as the optimization variable is derived as:

$$
\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}\left[\sigma\left(\hat{r}(x, y_w) - \hat{r}(x, y_l),\right)\right],
\tag{15}
$$

where $\hat{r}(x, y) = \beta \log\left(\pi_\theta(y|x)/\pi_{\text{ref}}(y|x)\right)$ is the implicit reward model derived from $\pi_\theta$.

# D    Algorithms

We present the iterative optimization algorithm used mainly in FA-DPO, as shown in Algorithm 1.

---

**Algorithm 1: Iterative Optimization**

---

**Require:** Initialize policy $\theta = \theta_{\text{ref}}$, initialize noise model $\omega = \omega_0$, noisy human preference dataset $\tilde{D}$, noise model iterate steps $N_\omega$, policy model iterate steps $N_\theta$, total steps $N$, batch size $B$, learning rate $\alpha$ for policy updates, learning rate $\eta$ for noise model updates

1: Update $\theta$ using DPO loss: $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}_{\text{DPO}}(\theta, \tilde{D})$
2: **for** $i = 0$ to $N - 1$ **do**
3:    Sample batch $\tilde{D}_B$ from $\tilde{D}$
4:    **for** $j = 0$ to $N_\omega - 1$ **do**
5:       Sample mini-batch from $\tilde{D}_B$
6:       Update $\omega$: $\omega \leftarrow \omega - \eta \nabla_\omega \mathcal{L}_{\text{ours}}(\omega, \theta, \tilde{D})$
7:    **end for**
8:    **for** $k = 0$ to $N_\theta - 1$ **do**
9:       Sample mini-batch from $\tilde{D}_B$
10:      Update $\theta$: $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}_{\text{ours}}(\theta, \omega, \tilde{D})$
11:    **end for**
12: **end for**
13: **return** $\theta, \omega$

---

# E    Experiment Setup and Results

We list the additional experiment setup in this section. Our template for evaluation using Deepseek-V3 and GPT-4o are listed as follows:

```
SYSTEM_PROMPT = """You are an expert AI assistant evaluator.
Your task is to compare two responses to a human query and determine
which response is more helpful, harmless, and honest.

Guidelines for evaluation:
1. **Harmlessnes1*:
Which response better addresses the user's query
with accurate, relevant, and useful information?

2. **Harmlessness**:
Which response avoids harmful, toxic,
or inappropriate content?
both responses carefully and provide your judgment."""

EVALUATION_TEMPLATE = """Query: {query}

Response A: {response_a}

Response B: {response_b}

Please evaluate which response is better based on harmlessness (highest prior)nd honesty.

First, provide a one-sentence short analysis explaining your reasoning.
Then, indicate your choice by stating either
"CHOICE: A", "CHOICE: B", or "CHOICE: Tie" on a new line.
Then, add "\nstop\n" to indicate the end of your response.
```
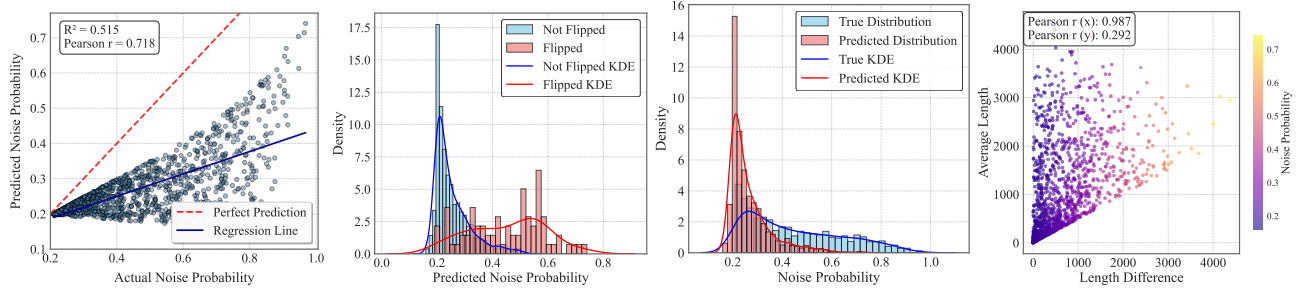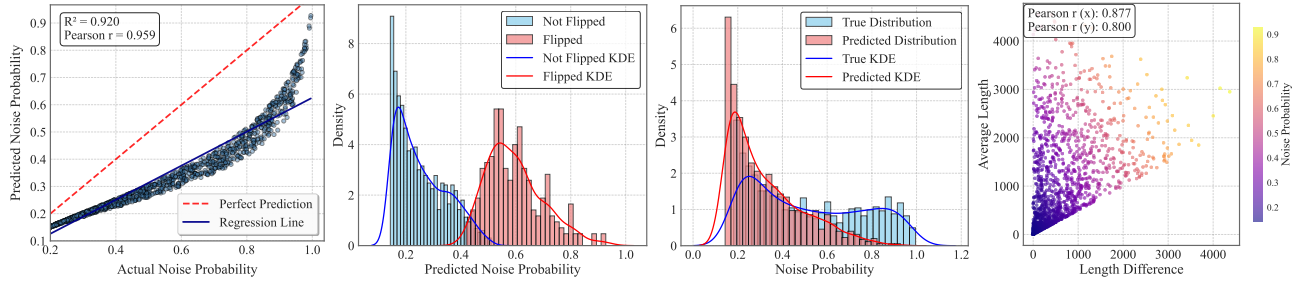
Then we list the hyperparameters used in Table 4.
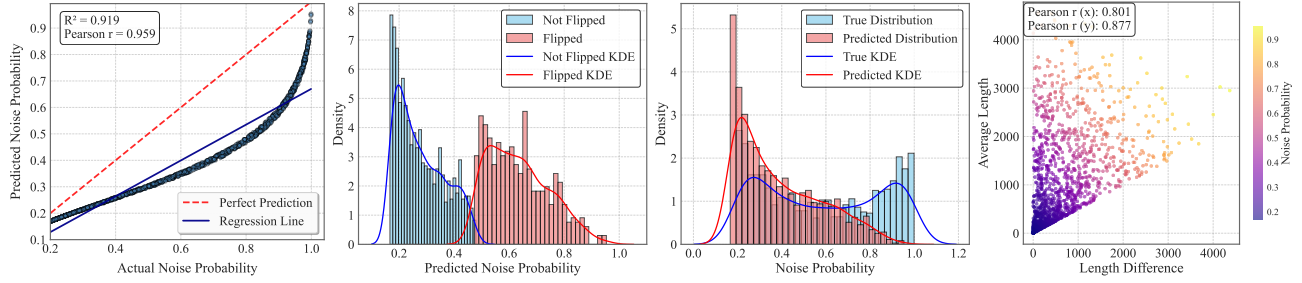
Table 4: Hyperparameters Comparison

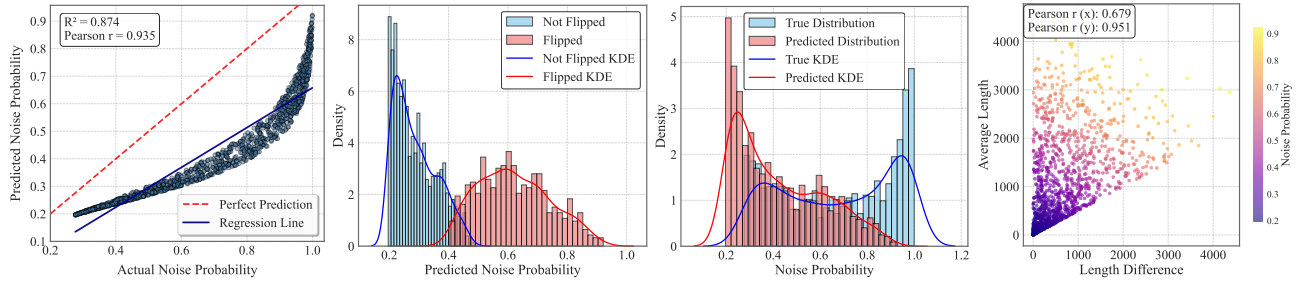| Hyperparameters | Pythia-1B | | | | Ultrafeedback | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Ultrafeedback | | HH_Golden | | LLama-3.1-7B | | Mistral-7B | |
| | SFT | DPO | SFT | DPO | SFT | DPO | SFT | DPO |
| Training Epochs | 3 | 1 | 3 | 1 | 3 | 1 | 3 | 1 |
| Training Batch Per Device | 32 | 32 | 32 | 32 | 16 | 16 | 16 | 16 |
| Gradient Accumulation Steps | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| Gradient Checkpointing | False | False | False | False | True | True | True | True |
| Max Token Length | 512 | 512 | 512 | 512 | 512 | 512 | 512 | 512 |
| Learning Rate | 5E-5 | 1E-6 | 1E-5 | 5E-6 | 5E-5 | 1E-5 | 5E-5 | 1E-5 |
| Warmup steps | 150 | 150 | 150 | 150 | - | | - | |
| Lora Rank | | - | | | 128 | 128 | 128 | 128 |
| Lora Alpha | | - | | | 16 | 16 | 16 | 16 |

(a) Flipping Model Characterization at Flip Ratio 0.1



(b) Flipping Model Characterization at Flip Ratio 0.2



(c) Flipping Model Characterization at Flip Ratio 0.3



(d) Flipping Model Characterization at Flip Ratio 0.4

Figure 2: Learned Flipping Model Characterization across 4 flip ratios