000 001

002 003

# **Uncertainty-aware Preference Alignment for Diffusion Policies**

Anonymous Authors<sup>1</sup>

#### Abstract

Recent advancements in diffusion policies have demonstrated promising performance in decisionmaking tasks. To align these policies with human preferences, a common approach is incorporating Preference-based Reinforcement Learning (PbRL) into policy tuning. However, since preference data is practically collected from populations with different backgrounds, a key challenge lies in handling the inherent uncertainties in people's preferences during policy updates. To address this challenge, we propose the Diff-UAPA algorithm, designed for uncertainty-aware preference alignment in diffusion policies. Specifically, Diff-UAPA introduces a novel iterative preference alignment framework in which the diffusion policy adapts incrementally to preferences from different user groups. To accommodate this online learning paradigm, Diff-UAPA employs a maximum posterior objective, which aligns the diffusion policy with regret-based preferences under the guidance of an informative Beta prior. This approach enables direct optimization of the diffusion policy without specifying any reward functions, while effectively mitigating the influence of inconsistent preferences across different user groups. We conduct extensive experiments across various robot control tasks and diverse human preference configurations, demonstrating the robustness and reliability of Diff-UAPA in achieving effective preference alignment.

# 1. Introduction

Reinforcement Learning (RL) algorithms commonly employ either deterministic or Gaussian policies to tackle sequential decision-making tasks by optimizing cumulative rewards (Sutton & Barto, 2018; Wang et al., 2022). Although these RL policies have demonstrated notable success across a wide range of applications (Mnih et al., 2015; Silver et al., 2016; Fang et al., 2019), they may struggle with learning multi-modal policies, which may hinder their ability to generalize effectively and lead to suboptimal performance in complex environments (Zhu et al., 2023). Recently, diffusion models have gained attention due to their strong modeling capabilities (Ho et al., 2020; Song et al., 2020). As a result, more studies have investigated the application of diffusion models in RL tasks, particularly in leveraging diffusion models as policies to model complex action distributions and behaviors (Wang et al., 2023; Chen et al., 2023a; Kang et al., 2023a; Lu et al., 2023; Chi et al., 2023). To learn a diffusion policy that generates desired outputs, recent approaches have leveraged Preference-based Reinforcement Learning (PbRL) (Christiano et al., 2017) techniques, which address a learning-to-rank problem using preference data, enabling alignment with human intentions (Wallace et al., 2024; Dong et al., 2024; Shan et al., 2024).

In practice, preferences are typically gathered from a diverse population, encompassing a wide range of expertise, perspectives, and beliefs. This diversity presents a significant challenge, as preferences from different user groups may conflict or evolve over time, introducing great uncertainties during policy updates. To ensure more reliable preference alignment, this necessitates the development of a policy that could account for the uncertainty arising from potentially inconsistent preferences. However, common PbRL approaches are typically based on the Bradley-Terry model (Bradley & Terry, 1952) with maximum likelihood estimation, which lacks sensitivity to the inherent uncertainties from preference datasets.

To address the uncertainties in preference alignment, several methods (Liang et al., 2022; Shin et al., 2023; Xue et al., 2024) have employed techniques such as ensemble models and Bayesian dropout. However, the underlying mechanism by which the estimated ensembles correlate with uncertainty remains largely unexplained. Motivated by recent work (Xu et al., 2025), which proposes learning a distributional reward model using a Maximum A Posteriori (MAP) objective to address epistemic uncertainty from an offline preference dataset, we explore how to bypass the reward learning and develop an uncertainty-aware algorithm beyond the offline setting for aligning diffusion policies.

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.



Figure 1. The framework of Diff-UAPA. Given the potentially inconsistent preference dataset ranked by diverse humans, we first learn a Beta prior to capture uncertainties, and then derive a Maximum A Posteriori (MAP) objective to align the diffusion policies.

065 In this work, we introduce Uncertainty-aware Preference 066 Alignment for Diffusion Policies (Diff-UAPA), a novel al-067 gorithm designed to align diffusion policies with human 068 preferences using an uncertainty-aware objective, as illus-069 trated in Figure 1. Specifically, we introduce an iterative 070 preference alignment framework, in which the diffusion policy progressively adapts to the labels coming from different user groups, each of which may have distinct preferences. To address this challenge, Diff-UAPA involves learning an 074 informative Beta prior that captures the uncertainty arising 075 from diverse human preferences. By interpreting prefer-076 ence alignment as a voting process, we demonstrate that 077 the Beta distribution is sensitive to the uncertainty among 078 compared trajectories, assigning high confidence to trajecto-079 ries in which the majority of human raters share a common preference and low confidence to those with divergent pref-081 erences. To ensure computational tractability, we parame-082 terize the Beta distribution with neural networks and train 083 the model via variational inference.

085 Guided by the informative Beta prior, Diff-UAPA aligns 086 the diffusion policy with a regret-based preference model, which inherently defines a unified Maximum A Posteriori 087 (MAP) objective. This method enables direct optimization 088 089 of the diffusion policy without requiring a reward function, while also effectively accounting for the uncertainties aris-090 091 ing from noisy preferences across diverse user groups.

092 To evaluate the empirical performance of Diff-UAPA, we 093 conduct extensive experiments across a diverse range of 094 robot manipulation and locomotion tasks, comparing its 095 performance against recently proposed baseline methods. 096 Furthermore, we investigate its effectiveness using hetero-097 geneous human preference data, including synthesized, re-098 alistic, and noisy preferences. The results demonstrate the 099 robustness and reliability of Diff-UAPA in handling varying 100 levels of uncertainty in preference data.

## 2. Related Works

109

061 062

063

064

#### 2.1. Preference-based Reinforcement Learning 104

Preference-based Reinforcement Learning (PbRL) is a piv-105 otal approach for aligning agents with human intent, particu-106 larly in scenarios where specifying explicit reward functions is challenging (MacGlashan et al., 2017; Warnell et al., 108

2018; Wirth et al., 2017). Previous works generally adopt a two-step procedure, where an explicit reward model is first inferred from human preferences using the Bradley-Terry model (Bradley & Terry, 1952), followed by training an RL agent to optimize the learned reward (Christiano et al., 2017; Ibarz et al., 2018). Building on this framework, several methods (Lee et al., 2021; Park et al., 2022; Hejna III & Sadigh, 2023; Liu et al., 2022; Liang et al., 2022; Hwang et al., 2023; Choi et al., 2024) have enhanced the learning process, focusing on improving efficiency and capability. In terms of preference modeling, while earlier works generally assume that preferences are generated based on the sum of Markovian rewards, recent studies (Kim et al., 2023; Verma & Metcalf, 2024) have proposed modeling preferences using non-Markovian rewards. Instead of learning an explicit reward model, another line of research focuses on directly optimizing policies or extracting value functions from human preferences (An et al., 2023; Hejna et al., 2024; Hejna & Sadigh, 2024). This approach is more straightforward, avoiding the biases and information bottleneck from intermediate reward modeling (Kang et al., 2023b).

#### 2.2. Diffusion Policy for Decision Making

Diffusion models have outperformed earlier generative models in both sample quality and training stability, gaining significant attention across various domains, including offline RL (Janner et al., 2022; Ajay et al., 2023), online RL (Yang et al., 2023; Chen et al., 2024), and robotics (Sridhar et al., 2024; Chen et al., 2023b; Xu et al., 2023). Recent advancements have leveraged diffusion models as RL policies to capture arbitrary action distributions, improving decision-making capabilities (Zhu et al., 2023). Among these works, Diffusion-QL (Wang et al., 2023), first integrated diffusion policies into the Q-learning framework. Following this, SfBC (Chen et al., 2023a) refined policy learning by decoupling behavior learning from action evaluation, while CEP (Lu et al., 2023) extended this framework to enable sampling from broader energy-guided distributions. CPQL (Chen et al., 2024) introduced consistency models to accelerate training and sampling, and EQP (Kang et al., 2023a) enhanced training efficiency with single-step model predictions for action approximations. In preferencebased tasks, AlignDiff (Dong et al., 2024) utilized diffusion

planners to generate trajectories aligned with human preferences through a two-step procedure, while FKPD (Shan
et al., 2024) introduced a one-step framework for direct
alignment. However, these methods often fail to account
for the uncertainties inherent in human preferences. How to
handle these uncertainties when aligning diffusion policies
remains a critical challenge (Casper et al., 2023).

# **3. Problem Formulation**

117

120 Preference-based Reinforcement Learning (PbRL). Re-121 inforcement Learning (RL) algorithms (Sutton & Barto, 122 2018) typically consider an episodic Markov Decision Pro-123 cess (MDP), which is formally defined as a tuple  $\mathcal{M}$  = 124  $(\mathcal{S}, \mathcal{A}, p_{\mathcal{R}}, p_{\mathcal{T}}, \gamma, T, \mu_0)$ , where: 1)  $\mathcal{S}$  and  $\mathcal{A}$  represent the 125 state and action spaces, 2)  $p_{\mathcal{R}}(r|s, a)$  and  $p_{\mathcal{T}}(s'|s, a)$  define 126 the (stochastic) reward and transition functions, 3)  $\gamma \in (0, 1]$ 127 is the discount factor, 4)  $\mu_0$  denotes the initial state distri-128 bution and 5)  $T \in (0,\infty)$  denotes a non-fixed planning 129 horizon, and the games is reset when the agent reaches a ter-130 minating or goal state at a time step T. In many applications, 131 the reward function is not directly available, reducing the 132 episodic MDP to a reward-free MDP  $\mathcal{M}_{/r}$ . To resolve this 133 challenge, PbRL algorithms (Christiano et al., 2017) pro-134 posed learning the reward function from human preferences 135 datatset. Specifically, given an unlabeled dataset of trajec-136 tory segments  $\mathcal{D}_{\tau} = \{\tau\}$ , humans randomly select a pair 137 of trajectories and rank them according to their preferences 138 on the optimality. By recording these pair-wise compar-139 isons, we create a preference dataset  $\mathcal{D}_{pref} = \{(\tau^w, \tau^l)\},\$ where each trajectory segment of length k is defined as 140 141  $\tau = (s_1, a_1, s_2, a_2, \dots, s_k, a_k)$ , and  $\tau^w$  is preferred over 142  $\tau^{l}$ . Based on this dataset, recent methods (Christiano et al., 143 2017; Ibarz et al., 2018) commonly infer the rewards by em-144 ploying the Bradley-Terry model (Bradley & Terry, 1952) 145 with maximum likelihood estimation (MLE). 146

147 Uncertainty Model in Preference Alignment. The 148 Bradley-Terry model (Bradley & Terry, 1952) can effec-149 tively model pairwise comparisons, whether by explicitly 150 inferring a reward function (Christiano et al., 2017; Lee 151 et al., 2021; Park et al., 2022) or by directly aligning poli-152 cies with preferences (Hejna et al., 2024; An et al., 2023). 153 However, this approach fails to account for the inherent 154 uncertainty in human preferences (Newman, 2023; Xu et al., 155 2025), particularly when these preferences are collected 156 from a diverse population with varying levels of expertise, 157 perspectives, and beliefs. More critically, for continuous learning, the policy must adapt dynamically to preferences 158 159 from different user groups, which often arrive incrementally 160 over time. To resolve these challenges, we study an iterative 161 preference alignment problem:

162 163 164 **Definition 3.1.** (Iterative Preference Alignment) Let  $\mathcal{D}_{\tau} = \tau$  $\tau$  denote the trajectory dataset, and let  $\mathcal{D}_{pair}^n = (\tau^i, \tau^j)$  represent the pairwise comparisons dataset constructed at the  $n^{th}$  iteration. These comparisons are generated by 1) sampling pairs of trajectories from  $D\tau$  and 2) inviting a group of annotators to label them. The algorithm must progressively align the policy  $\pi$  with the preference dataset  $D_{\text{pair}}^n$  at each round  $n \in [1, N]$  in an online manner.

In this setting, different groups of human annotators may provide inconsistent or even conflicting preferences for the same pair of trajectories (Liang et al., 2022; Shin et al., 2023; Xue et al., 2024). The problem solver must dynamically adapt the policy to iteratively updated preference signals while ensuring that the learned policy effectively represents general human preferences by performing online updates.

Additionally, apart from the preference signals, the trajectory dataset  $D_{\tau}$  can in principle be updated based on interaction from the environment. However, in practice, such interactions are not always available, and thus we assume  $D_{\tau}$  mainly records only offline trajectories. The primary challenge is to stabilize the policy optimization process and learn a reliable control policy by effectively managing the aleatoric uncertainty inherent in stochastic and potentially inconsistent preference signals on the provided trajectories.

**Preference Alignment for Diffusion Policies.** While previous PbRL methods have commonly focused on policies modeled by feed-forward neural networks, recent studies highlight the superior control performance achieved by diffusion-based policies (Zhu et al., 2023). Denoising diffusion models (Ho et al., 2020) represent a class of generative models characterized by an iterative diffusion and denoising process. Diffusion models have gained significant attention in decision-making tasks due to their ability to represent complex multi-modal distributions (Zhu et al., 2023). This capability is crucial for characterizing the policy function  $\pi_{\theta}(a|s)$ , surpassing previous deterministic or Gaussian-based policies (Chi et al., 2023; Wang et al., 2023). Diffusion policies are typically formulated as conditional generative models as follows<sup>1</sup>:

$$\pi_{\theta}(a_t|s_t) = \int \mathcal{N}(a_t^I; \mathbf{0}, \mathbf{I}) \prod_{i=1}^{I} \pi_{\theta}(a_t^{i-1}|a_t^i, s_t) \mathrm{d}a_t^{1:I},$$
<sup>(1)</sup>

where  $\pi_{\theta}(a_t^{i-1}|a_t^i, s_t)$  is often parameterized as Gaussian with fixed timestep-dependent covariances as  $\mathcal{N}(a_t^{i-1}|\mu_{\theta}(a_t^i, s_t, i), \Sigma^i)$ . Although diffusion policies can be trained from offline datasets, their performance is often constrained by the size, quality, and availability of the expert demonstration dataset. As a result, many previous methods have utilized RL algorithms to improve these policies with experience data sampled from an interactive MDP

<sup>&</sup>lt;sup>1</sup>In this work, we use superscripts  $(i \in \{0, 1, ..., I\}$  to denote diffusion timesteps and subscripts  $(t \in \{0, 1, ..., T\})$  to denote trajectory timesteps.

165 environment (Kang et al., 2023a; Psenka et al., 2024). In this setting, recent research (Wallace et al., 2024) proposed 167 leveraging Direct Preference Optimization (DPO) (Rafailov 168 et al., 2023) to align diffusion policies with human prefer-169 ences based on  $\mathcal{D}_{pref}$ . Specifically, DPO algorithms directly 170 optimize policies without learning a reward model, thereby 171 significantly enhancing the efficiency and stability of the 172 training process. To train  $\pi_{\theta}$ , the maximum likelihood ob-173 jective for state-action pairs is defined as follows:

$$\begin{array}{ll}
 174 \\
 175 \\
 176 \\
 177 \\
 177 \\
 178 \\
 179 \\
 179 \\
 179
\end{array}$$

$$\begin{array}{ll}
 L(\theta) = -\mathbb{E}\Big[\log\sigma\Big(-\lambda I \cdot (2) \\
 (\|\epsilon^w - \epsilon_\theta(a^{i,w}, s^w, i)\|_2^2 - \|\epsilon^w - \epsilon_{\rm ref}(a^{i,w}, s^w, i)\|_2^2) \\
 128 \\
 - (\|\epsilon^l - \epsilon_\theta(a^{i,l}, s^l, i)\|_2^2 - \|\epsilon^l - \epsilon_{\rm ref}(a^{i,l}, s^l, i)\|_2^2) \Big)\Big],$$

$$(2)$$

180 where 1)  $((s^w, a^{0,w}), (s^l, a^{0,l})) \sim \mathcal{D}_{\text{pref}}$  are state-action samples from preference dataset, 2)  $i \sim \mathcal{U}(0, I)$  is the 182 diffusion timestep, and 3)  $a^{i,w/l} \sim q(a^{i,w/l}|a^{0,w/l},s^{w/l})$ 183 denotes the action  $a^{0,w/l}$  corrupted with noise  $\epsilon^{w/l}$  after i 184 diffusion steps, as defined in (Ho et al., 2020). In this study, 185 we explore addressing the iterative preference alignment 186 problem by aligning human preferences with a diffusion 187 policy model. 188

181

189

190

191

213 214 215

216

217

218

219

# 4. Uncertainty-Aware Preference Alignment for Diffusion Policies

192 In this section, we outline our approach for aligning a dif-193 fusion policy with human preferences while effectively ac-194 counting for uncertainty. Specifically, we present: 1) a Max-195 imum Likelihood Estimation (MLE) objective for diffusion 196 policy alignment, based on maximum entropy framework 197 and direct preference optimization (Section 4.1), 2) a Maxi-198 mum A Posteriori (MAP) objective that incorporates a Beta 199 prior model for capturing the underlying uncertainties (Sec-200 tion 4.2), and 3) the training procedure for the Beta prior 201 model (Section 4.3). 202

#### 203 4.1. Maximum Likelihood Diffusion Policy Alignment

204 MaxEnt Alignment under Regret Preference. Follow-205 ing previous works on preference alignment (Hejna et al., 206 2024; Rafailov et al., 2024; Ouyang et al., 2022), we adopt 207 the Maximum Entropy (MaxEnt) RL framework. In this 208 approach, the objective is to learn a policy  $\pi_{\theta}$  that not only 209 maximizes its cumulative discounted rewards but also in-210 corporates the causal entropy, while regularizing the KL-211 divergence from a reference policy (Ziebart, 2010): 212

$$\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{T} \gamma^t (r(s_t, a_t) - \alpha \log \frac{\pi(a_t|s_t)}{\pi_{\text{ref}}(a_t|s_t)}) \right], \quad (3)$$

Here,  $\alpha$  determines the weight of entropy in the optimization objective. Upon learning an optimal policy  $\pi^*$ , we can compute the corresponding optimal state-value function  $V^*(s_t)$ , the optimal state-action value function  $Q^*(s_t, a_t)$ , and the

optimal advantage function  $A^*(s_t, a_t) \triangleq Q^*(s_t, a_t) V^*(s_t)$ . More importantly, in the MaxEnt RL setting, the optimal advantage function is proportional to the loglikelihood of the optimal and reference policy (Haarnoja et al., 2017; Hejna et al., 2024):

$$A^{*}(s_{t}, a_{t}) = \alpha \log \frac{\pi^{*}(a_{t}|s_{t})}{\pi_{\text{ref}}(a_{t}|s_{t})}.$$
(4)

To stabilize the process of preference alignment, we follow (Knox et al., 2022) and base the preference alignment on discounted regrets, defined as  $-\sum \gamma^t (V(s_t) - Q(s_t, a_t))$ . In this framework, a trajectory segment is preferred if it incurs lower regret compared to the intended optimal policy, so that the preference between trajectory segments  $(\tau^w, \tau^l)$ can be modeled as:

$$P_{A*}(\tau^{w} \succ \tau^{l}) =$$

$$\frac{\exp \sum_{t=0}^{k} \gamma^{t} A^{*}(s_{t}^{w}, a_{t}^{w})}{\exp \sum_{t=0}^{k} \gamma^{t} A^{*}(s_{t}^{w}, a_{t}^{w}) + \exp \sum_{t=0}^{k} \gamma^{t} A^{*}(s_{t}^{l}, a_{t}^{l})}.$$
(5)

By substituting Equation (4) into Equation (5), the advantage function  $A^*$  can be replaced by the optimal policy  $\pi^*$  under the MaxEnt framework. The learned policy  $\pi_{\theta}$ can then be optimized through maximum the likelihood of generating preferences as follows (Hejna et al., 2024):

$$\mathcal{L}_{CPL}^{(\tau^w,\tau^l)}(\theta) = -\log\sigma\Big(\alpha \cdot (6) \\ \left(\sum_{t=0}^k \gamma^t \log \frac{\pi_{\theta}(a_t^w | s_t^w)}{\pi_{ref}(a_t^w | s_t^w)} - \sum_{t=0}^k \gamma^t \log \frac{\pi_{\theta}(a_t^l | s_t^l)}{\pi_{ref}(a_t^l | s_t^l)}\right)\Big)$$

**Diffusion Policy Alignment.** To adapt the previous model to aligning the diffusion policy  $\pi_{\theta}(a_t|s_t)$  as defined in Equation (1), a primary difficulty is due to the intractability of diffusion policy  $\pi_{\theta}(a_t|s_t) = \int \pi_{\theta}(a_t^{0:I}|s_t) da_t^{1:I}$ , as it requires marginalizing over all possible diffusion paths  $(a_t^1, a_t^2, \dots, a_t^I)$  that lead to  $a_t^0$ . To address it, we propose modeling the chain reward function (Wallace et al., 2024):

$$r(s_t, a_t^0) = \mathbb{E}_{\pi_\theta(a_t^{1:I} | a_t^0, s_t)}[r(s_t, a_t^{0:I})].$$
(7)

The optimal chain advantage function can be defined as:

$$A^{*}(s_{t}, a_{t}^{0}) = \mathbb{E}_{\pi^{*}_{\theta}(a_{t}^{1:I}|a_{t}^{0}, s_{t})} \left[ A^{*}(s_{t}, a_{t}^{0:I}) \right]$$
(8)

$$= \mathbb{E}_{\pi_{\theta}^{*}(a_{t}^{1:I}|a_{t}^{0},s_{t})} \left[ \alpha \log \frac{\pi_{\theta}^{*}(a_{t}^{0:I}|s_{t}))}{\pi_{\text{ref}}(a_{t}^{0:I}|s_{t}))} \right].$$
(9)

In principle, we can interpret the latent diffusion actions as a unified chain action  $\overline{a_t} = a_t^{0:I}$ , despite the final output being determined by  $a_t^0$ . This perspective allows us to reformulate Equation (3) in terms of the diffusion policy:

$$\max_{\pi_{\theta}} \mathbb{E}_{\pi_{\theta}(\overline{a_t}|s_t)} \left[ \sum_{t=0}^{T} \gamma^t (r(s_t, \overline{a_t}) - \alpha \log \frac{\pi_{\theta}(\overline{a_t}|s_t))}{\pi_{\text{ref}}(\overline{a_t}|s_t))} \right].$$
(10)

This objective is defined over the entire diffusion path  $\overline{a_t}$ , which aims to maximize the cumulative rewards and the entropy within a trajectory across the reverse process.

220 By paralleling from Equation (3) to Equation (6), the objective in (10) can be directly optimized with respect to the diffusion policy  $\pi_{\theta}(\overline{a_t}|s_t)$  by maximizing the following likelihood:

$$\mathcal{L}_{1,\text{MLE}}^{(\tau^{w},\tau^{l})}(\theta) = -\log\sigma\Big(\alpha \cdot \tag{11})$$

$$\Big(\sum_{t=0}^{k} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,w}|s_{t}^{w},a_{t}^{0,w})} \left[\gamma^{t}\log\frac{\pi_{\theta}(\overline{a_{t}^{w}}|s_{t}^{w})}{\pi_{\text{ref}}(\overline{a_{t}^{w}}|s_{t}^{w})}\right]$$

$$-\sum_{t=0}^{k} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,l}|s_{t}^{l},a_{t}^{0,l})} \left[\gamma^{t}\log\frac{\pi_{\theta}(\overline{a_{t}^{l}}|s_{t}^{l})}{\pi_{\text{ref}}(\overline{a_{t}^{l}}|s_{t}^{l})}\right]\Big)\Big),$$

 $\overline{t=0}$ 

where  $\sigma$  is the sigmoid function. However, major challenges in optimizing this objective lie in: 1) inefficiency, due to the sequential computation required across many timesteps, and 2) intractability, stemming from the need to evaluate the joint distribution. Inspired by Wallace et al. (2024), we leverage Jensen's inequality and the convexity of the  $-\log \sigma$  function to move the expectation operator outside, thereby improving efficiency. Additionally, we approximate the reverse process  $\pi_{\theta}(a_t^{1:I}|s_t)$  using the forward process  $q(a_t^{1:I}|s_t)$ , which makes the problem more tractable. With some algebra, we derive the following loss function:

$$\mathcal{L}_{1,\text{MLE}}^{(\tau^{w},\tau^{l})}(\theta) \leq -\mathbb{E}_{a_{t}^{i,w} \sim q(a_{t}^{i,w}|a_{t}^{0,w},s_{t}^{w}),} \left[ \log \sigma \left( -\alpha I \cdot \frac{1}{a_{t}^{i,l} \sim q(a_{t}^{i,l}|a_{t}^{0,l},s_{t}^{l})} \right) \left( \sum_{t=0}^{k} \gamma^{t} (\|\epsilon^{w} - \epsilon_{\theta}(a_{t}^{i,w},s_{t}^{w},i)\|_{2}^{2} - \|\epsilon^{w} - \epsilon_{\text{ref}}(a_{t}^{n,w},s_{t}^{w},i)\|_{2}^{2}) - \sum_{t=0}^{k} \gamma^{t} (\|\epsilon^{l} - \epsilon_{\theta}(a_{t}^{i,l},s_{t}^{l},i)\|_{2}^{2} - \|\epsilon^{l} - \epsilon_{\text{ref}}(a_{t}^{i,l},s_{t}^{l},i)\|_{2}^{2}) \right) \right) = \mathcal{L}_{2,\text{MLE}}^{(\tau^{w},\tau^{l})}(\theta),$$
(12)

The detailed deviation is shown in Appendix A.

#### 4.2. Bayesian Alignment with Informative Beta Prior

The regret preference model (Equation (5)) represents the likelihood of generating human preferences based on the advantage function. The corresponding maximum likelihood objective implicitly assumes a uniform prior over  $\sum_{t=0}^{k} \gamma^{t} A^{*}(s_{t}, a_{t})$ , which does not account for the uncertainty within the preference dataset, and may lead to divergence in the parameters of the learned policy (Newman, 2023; Xu et al., 2025). We present how to derive a more informative prior as follows.

Since human feedback is based on two trajectories rather than individual state-action pairs, we assume that the strength of a trajectory is defined by its trajectory-level advantage, represented by its discounted cumulative advantages under the diffusion policy  $\pi_{\theta}$ :

$$A^{\pi_{\theta}}(\tau) = \sum_{t=0}^{k} \gamma^{t} A^{\pi_{\theta}}(s_{t}, a_{t})$$
$$= \sum_{t=0}^{k} \gamma^{t} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I} | a_{t}^{0}, s_{t})} \left[ A^{\pi_{\theta}}(s_{t}, \overline{a_{t}}) \right].$$
(13)

The average strength of the trajectories under policy  $\pi_{\theta}$  is then defined as:

$$\bar{A}^{\pi_{\theta}} = \mathbb{E}_{\tau \sim \mathcal{D}_{\tau}} A_{\theta}(\tau) = \frac{1}{|\mathcal{D}_{\text{pref}}|} \sum_{\tau \in \mathcal{D}_{\text{pref}}} A^{\pi_{\theta}}(\tau).$$
(14)

Therefore, the probability of a trajectory with strength  $A^{\pi_{\theta}}(\tau)$  winning against the average candidate is  $\phi(\tau) =$  $\sigma(A^{\pi_{\theta}}(\tau) - \bar{A}^{\pi_{\theta}}) \in (0, 1)$ . By applying the chain rule, the prior on the advantage function can be defined as:

$$p_0(A^{\pi_{\theta}}(\tau)) = p_0(\phi(\tau)) \frac{\mathrm{d}\phi(\tau)}{\mathrm{d}A^{\pi_{\theta}}(\tau)}$$
$$= p_0(\phi(\tau))\sigma'(A^{\pi_{\theta}}(\tau) - \bar{A}^{\pi_{\theta}})(1 - \frac{1}{|\mathcal{D}_{\mathrm{pref}}|}).$$
(15)

This prior reflects our initial belief about the strength of different trajectories within the dataset. Motivated by Xu et al. (2025), we use the Beta distribution as the informative prior, i.e.,  $p_0(\phi(\tau)) = \text{Beta}(\phi(\tau); \alpha, \beta)$ . The main benefits of the Beta distribution are: 1) it is the conjugate prior for the Bernoulli distribution, and  $\phi(\tau)$  naturally ranges from (0, 1), which simplifies updates with new evidence, and 2) the parameters  $\alpha$  and  $\beta$  can intuitively represent the counts of preferred and unpreferred human feedback. By reformulating Eq. (15), we present the following proposition:

**Proposition 4.1.** Let the informative prior  $p_0(\phi(\tau))$  be a Beta distribution  $Beta(\phi(\tau); \alpha, \beta)$ . This prior can effectively capture the uncertainty arising from the iterative preference alignment process (Definition 3.1). Consequently, the prior on the strength of a trajectory is proportional to Beta( $(\phi(\tau); \alpha+1, \beta+1)$ ), i.e.,  $p_0(A^{\pi_{\theta}}(\tau)) \propto$  $Beta(\phi(\tau); \alpha + 1, \beta + 1).$ 

The proof is shown in Appendix C. The corresponding prior loss can then be derived in a manner similar to the derivation of the maximum likelihood loss (Eq. 11):

$$\mathcal{L}_{1,\text{prior}}^{\tau}(\theta) = -\log \operatorname{Beta}(\phi(\tau); \alpha + 1, \beta + 1) \\ \leq -\mathbb{E} \bigg[ \log \operatorname{Beta} \bigg( \sigma \bigg( -\alpha I \cdot \bigg( \sum_{t=0}^{k} \gamma^{t} (\|\epsilon - \epsilon_{\theta}(a_{t}^{i}, s_{t}, i)\|_{2}^{2} - \left\|\epsilon - \epsilon_{\text{ref}}(a_{t}^{i}, s_{t}, i)\|_{2}^{2} \bigg) - \sum_{\tau \in \mathcal{D}_{\text{pref}}, t=0}^{k} \frac{\gamma^{t}}{|\mathcal{D}_{\text{pref}}|} (\|\epsilon - \epsilon_{\theta}(a_{t}^{i}, s_{t}, i)\|_{2}^{2} - \left\|\epsilon - \epsilon_{\text{ref}}(a_{t}^{i}, s_{t}, i)\|_{2}^{2} \bigg) \bigg) \bigg); \alpha + 1, \beta + 1 \bigg) \bigg] \\ = \mathcal{L}_{2,\text{prior}}^{\tau}(\pi_{\theta})$$
(16)

324

325

327

329

Appendix B shows the detailed proof. Equation (16) can be interpreted as guiding the policy to align the estimated advantage function for trajectories with their prior distribution. Since  $P_{\text{MAP}}(A(\tau)) \propto p_0(A(\tau)) \cdot P_{\text{MLE}}(A(\tau))$ , by incorporating the prior into the MLE objective and maximizing the log form of the posterior, we can derive the Diff-UAPA loss:

$$\mathcal{L}_{\text{Diff-UAPA}}(\theta) = \mathbb{E}_{(\tau^{w},\tau^{l})\sim\mathcal{D}_{\text{pref}}} \Big[ \mathcal{L}_{2,\text{MLE}}^{(\tau^{w},\tau^{l})}(\pi_{\theta}) \\ + \mathcal{L}_{2,\text{prior}}^{\tau^{w}}(\pi_{\theta}) + \mathcal{L}_{2,\text{prior}}^{\tau^{l}}(\pi_{\theta}) \Big].$$
(17)

Maximizing the posterior probability, rather than the likelihood, incorporates prior knowledge and regularizes advantage values, preventing divergence. We introduce how to estimate the Beta prior in the following section.

#### 4.3. Training the Beta Prior Model

To learn the Beta prior  $p_0(\phi(\tau)|\mathcal{D}_{pref}) = \text{Beta}(\phi(\tau); \alpha, \beta)$ in continuous spaces, following (Xu et al., 2025), we propose using a variational inference approach to approximate it by estimating the approximate posterior  $q_{\xi}(\phi(\tau)|\mathcal{D}_{pref})$ , i.e.,  $p_0(\phi(\tau)|\mathcal{D}_{\text{pref}}) \simeq q_{\xi}(\phi(\tau)|\mathcal{D}_{\text{pref}})$ , where  $\xi$  is the model parameters. The objective is to minimize the Kullback-Leibler (KL) divergence between the prior and posterior, which is equivalent to maximizing the Evidence Lower Bound (ELBO). This leads to the following interpretation of the corresponding trajectory-wise objective (Xu et al., 2025):

$$\max_{\xi} \mathbb{E}_{\tau} \left[ \mathbb{E}_{q_{\xi}, (\tau^w, \tau^l) \in \mathcal{D}_{\text{pref}}} [\log \phi(\tau^w)] - \right]$$
(18)

$$\mathbb{E}_{q_{\xi},(\tau^{w},\tau^{l})\in\mathcal{D}_{\text{pref}}}[\log\phi(\tau^{l})] - D_{\text{KL}}[q_{\xi}(\phi(\tau)|\tau) \parallel p(\phi(\tau))]],$$

where 1)  $q_{\xi}(\phi(\tau)|\tau) = \text{Beta}(\alpha_{\tau}, \beta_{\tau})$ , where  $[\alpha_{\tau}, \beta_{\tau}] = f_{\xi}^{\text{Beta}}(\tau)$  and  $f_{\xi}^{\text{Beta}}$  denotes a neural network, 2)  $p(\phi(\tau)) = f_{\xi}^{\text{Beta}}(\tau)$ Beta $(\alpha_0, \beta_0)$ , with  $\alpha_0, \beta_0$  specifying our prior belief (we set  $\alpha_0 = \beta_0 = 1$  in this work), and 3)  $\phi(\tau)$  represents the Bernoulli probability that  $\tau^w$  is ranked higher than  $\tau^l$ . The first two terms aim to optimize the parameter  $\xi$  to align with the preference dataset, while the final KL-divergence term ensures the posterior distribution does not deviate too far from the prior belief, which can be optimized using the Dirichlet VAE approach (Joo et al., 2020).

In this work, we implement  $f_{\xi}^{\text{Beta}}(\tau)$  using a transformer-based neural network (Vaswani, 2017), where the trajectory  $\tau$  is fed as input and  $[\alpha_{\tau}, \beta_{\tau}]$  is produced as the output to form the Beta prior distribution. The complete Diff-UAPA algorithm is shown in Algorithm 1.

#### **5.** Empirical Evaluation

In this section, we empirically evaluate the proposed Diff-UAPA algorithm on four robot manipulation tasks across two environments (Section 5.1) and locomotion tasks with

Algorithm 1 Uncertainty-aware Preference Alignment for **Diffusion Policies (Diff-UAPA)** 

- 1: **Input:** Trajectory dataset  $\mathcal{D}_{\tau}$ , preference dataset  $\mathcal{D}_{\text{pref}}$ , prior training epochs M, policy training epochs N. 2: Initialize Beta prior model  $f_{\xi}^{\text{Beta}}(\tau)$ , reference policy
- $\pi_{\text{ref}}(a|s)$ , and diffusion policy  $\pi_{\theta}(a|s)$ .
- 3: Learn  $\pi_{ref}$  based on  $\mathcal{D}_{\tau}$  through behavior cloning.
- 4: for  $m = 1, \cdots, M$  do
- Update the Beta prior  $f_{\xi}^{\text{Beta}}$  with objective (18). 5:
- 6: end for
- 7: for  $n = 1, \dots, N$  do
- Update the diffusion policy  $\pi_{\theta}$  by minimizing 8: Eq. (17).

9: end for

real human preferences (Section 5.2), where preferences are continuously updated and may exhibit inconsistencies. Additionally, we evaluate the noise sensitivity of the proposed method under different levels of preference inconsistency (Section 5.3).

**Experiment Settings.** We evaluate the methods on three tasks in Robomimic (Mandlekar et al., 2021) and one longhorizon Franka Kitchen (Gupta et al., 2019) environment for manipulation tasks, as well as two environments in D4RL (Fu et al., 2020) with real human preferences for locomotion tasks. Our experiments consist of four rounds of iterative updates, with each round consisting of a fixed number of training episodes. To account for potential inconsistencies in human preferences, we introduce a reverse rate into the ground-truth preference data. Specifically, in each update round, we randomly select 20% of trajectory pairs and apply a 50% reversal rate by swapping the winner and the loser. The learning rate is reset at the beginning of each round to enhance stability and convergence. After training, the policy is evaluated over 10 episodes in 56 parallel environments. Each experiment is repeated using three different random seeds, and the mean  $\pm$  standard deviation (std) of the results is reported. More experimental details can be found in Appendix D.1.

Comparison Methods. We utilize two baseline policies: the Gaussian-based policy from Behavior Transformer (BET) (Shafiullah et al., 2022) and the Diffusion Policy (Diff) (Chi et al., 2023). In BET, we apply focal loss (Mukhoti et al., 2020) for preference-based learning and leverage the full set of trajectories in the preference dataset for training the diffusion policy.

Building on BET, we propose the following comparison methods: 1) BET-Direct Preference Optimization (**BET-DPO**) and 2) BET-Contrastive Preference Learning (BET-CPL), which leverage direct preference optimization

Table 1. Success rates (in percentage) of all methods across the Robomimic and Kitchen tasks, with each value presented as the mean  $\pm$  std, computed over 3 training seeds and 560 evaluation episodes. The best results for each task are highlighted in bold. For the Kitchen task, px indicates the frequency of interaction with x or more objects.

| 334        |             |                | Robomimic      |                                  |                | Kito                             | chen                             |                                  |
|------------|-------------|----------------|----------------|----------------------------------|----------------|----------------------------------|----------------------------------|----------------------------------|
| 335<br>336 |             |                |                |                                  |                |                                  |                                  |                                  |
| 337        |             |                | / THAT         | Carl Star                        |                |                                  |                                  |                                  |
| 338        |             |                |                |                                  |                | 20536                            | 100                              |                                  |
| 339<br>340 |             | L :ft          | Con            | Squara                           |                | n)                               | n <sup>2</sup>                   |                                  |
| 341        |             | LIII           | Can            | Square                           | p1             | p2                               | p3                               | p4                               |
| 3/12       | BET         | $43.6 \pm 3.8$ | $48.8\pm3.1$   | $55.1\pm2.0$                     | $96.4 \pm 1.2$ | $96.2\pm1.0$                     | $76.6 \pm 1.3$                   | $44.6\pm2.0$                     |
| 242        | BET-CPL     | $49.2 \pm 4.4$ | $42.1\pm1.1$   | $57.6\pm2.3$                     | $97.0 \pm 1.0$ | $96.4\pm0.5$                     | $88.4\pm2.3$                     | $62.6\pm2.0$                     |
| 243        | BET-DPO     | $43.7 \pm 3.3$ | $47.0\pm1.0$   | $42.7\pm3.6$                     | $85.5 \pm 8.5$ | $84.8\pm8.7$                     | $80.9\pm9.4$                     | $57.4\pm6.6$                     |
| 344        | Diff        | $45.1 \pm 3.0$ | $47.9\pm2.3$   | $52.8\pm2.9$                     | $99.2 \pm 0.8$ | $98.4 \pm 1.1$                   | $91.8\pm0.8$                     | $59.0 \pm 1.1$                   |
| 345        | Diff-CPL    | $48.6 \pm 2.2$ | $45.9\pm2.8$   | $55.2\pm5.7$                     | $100.0\pm0.0$  | $99.6\pm0.2$                     | $94.2\pm0.2$                     | $63.5\pm0.8$                     |
| 346        | FKPD        | $51.2 \pm 0.7$ | $58.5\pm2.5$   | $64.4 \pm 2.7$                   | $99.8 \pm 0.3$ | $98.3 \pm 1.4$                   | $89.5\pm2.9$                     | $64.1 \pm 3.2$                   |
| 347        | Diff-UAPA-C | 56.1±0.9       | $61.3\pm2.2$   | $\textbf{68.1} \pm \textbf{0.6}$ | $100.0\pm0.0$  | $99.7\pm0.2$                     | $95.4\pm0.6$                     | $70.9\pm2.5$                     |
| 348        | Diff-UAPA-I | $54.3 \pm 1.1$ | $59.9 \pm 1.7$ | $66.2\pm1.3$                     | 99.9±0.1       | $\textbf{99.8} \pm \textbf{0.2}$ | $\textbf{95.7} \pm \textbf{1.9}$ | $\textbf{71.7} \pm \textbf{4.6}$ |
| .1 + 7     |             |                |                |                                  |                |                                  |                                  |                                  |

351 (Rafailov et al., 2023) and contrastive preference learning 352 (Hejna et al., 2024) to align the BET model. For diffusion-353 based policies, we introduce: 3) Diffusion Policy-CPL (Diff-354 CPL) that uses the MLE loss for aligning the diffusion 355 policy (Obj. 12), and 4) FKPD (Shan et al., 2024) that per-356 forms forward KL regularized preference optimization. For 357 our Diff-UAPA algorithm, we explore two distinct strate-358 gies for updating the Beta prior model: 5) Diff-UAPA-C 359 that trains the Beta model using full preference data across 360 the iterations without updates, and 6) Diff-UAPA-I that in-361 crementally updates the Beta model on the current noisy 362 preference data through the iterative process. 363

333

350

# <sup>364</sup> 5.1. Model Performance in Robot Manipulation Tasks

Task Description. In this experiment, we evaluate the model's performance across three tasks from 367 Robomimic (Mandlekar et al., 2021) and the Franka Kitchen task introduced in (Gupta et al., 2019), both of which use 369 state-based observations. Specifically, the three Robomimic 370 tasks-Lift, Can, and Square-address different manipu-371 lation challenges in a simulated environment, including object lifting, can manipulation, and square positioning. On 373 the other hand, the Franka Kitchen task involves complex, 374 multi-step, long-horizon activities that require interactions 375 with seven distinct objects, with the objective to complete 376 as many demonstrated tasks as possible, regardless of the execution order. Following Chi et al. (2023), we use suc-378 cess rate as the primary evaluation metric. For each task, 379 the reference policy  $\pi_{ref}$  is trained to achieve a success rate 380 of approximately 40%. We then roll out the policy to col-381 lect 560 trajectories per task and construct the preference 382 dataset based on their rewards. Please check Appendix D.2 383 384

for environmental details and Appendix D.3 for details on preference dataset construction.

**Results Analysis.** Table 1 presents the evaluation performance across three Robomimic tasks and the more complex Kitchen task. The results indicate that both variants of Diff-UAPA consistently outperform other methods across different tasks. This is primarily due to their use of a Beta prior, which effectively captures the uncertainty arising from potentially inconsistent preferences, thereby enhancing the diffusion policy training process. Moreover, the performance gap between Diff-UAPA-C and Diff-UAPA-I is relatively small, suggesting that the Beta prior can be trained effectively in both approaches, depending on the specific practice. This flexibility enhances the practical applicability of the proposed method. Notably, for the long-horizon Kitchen task, Diff-UAPA-I, which trains the Beta model incrementally, slightly outperforms Diff-UAPA-C, which pre-trains the Beta model using the complete dataset. This difference can be attributed to the fact that incremental training allows the model to adapt more dynamically to the changing preferences and environmental conditions over time, whereas pre-training may not fully capture such variability. We also provide the visualization results in Figure 2 in Appendix D.5

#### 5.2. Model Performance in Locomotion Tasks

**Task Description.** The primary goal of Preference-based Reinforcement Learning (PbRL) is to align policies with *human* preferences. In this section, we assess the performance of Diff-UAPA using real human preferences provided by the Uni-RLHF benchmark (Yuan et al., 2024) in the HalfCheetah and Walker environments from the D4RL

|             | BET            | BET-CPL      | BET-DPO     | Diff           | Diff-CPL     | FKPD         | Diff-UAPA-C                     | Diff-UAPA-I    |
|-------------|----------------|--------------|-------------|----------------|--------------|--------------|---------------------------------|----------------|
| HalfCheetah | $2577 \pm 198$ | $2976\pm 66$ | $2948\pm37$ | $2838\pm325$   | $3121\pm148$ | $3060\pm201$ | $\textbf{3399} \pm \textbf{72}$ | $3297 \pm 101$ |
| Hopper      | $1161\pm90$    | $1226\pm85$  | $1129\pm79$ | $1296 \pm 137$ | $1313\pm103$ | $1370\pm120$ | $1591\pm51$                     | $1499\pm70$    |

Table 2. Episodic rewards of all methods in the HalfCheetah and Hopper environments with real human preferences.

benchmark (Fu et al., 2020). To ensure the dataset encompasses a diverse range of trajectories for meaningful comparison, we use *medium-expert* datasets for both environments. These datasets combine expert demonstrations from a near-optimal policy with suboptimal data generated by a medium-performing policy. Please check Appendix D.2 for more environmental details.

**Results Analysis.** The empirical results for the locomotion tasks are presented in Table 2. We observe that Diff-UAPA consistently outperforms other baselines across both environments. The key reason for this is that, during the iterative preference alignment process, some trajectory pairs may receive inconsistent preference labels. These noisy labels introduce greater uncertainty, making it challenging for the policy to accurately assess the true value of these trajectories and replicate the higher-performing ones. Diff-UAPA effectively addresses this challenge by leveraging a prior model that captures this uncertainty, enabling the policy to evaluate the trajectories more fairly and reliably, which in turn leads to improved overall performance. We also observe that diffusion-based policies generally achieve better results than Gaussian-based policies, primarily due to their superior modeling capabilities, which becomes more crucial when accounting for underlying uncertainties.

#### 5.3. Experiments on Noise Sensitivity

**Task Description.** In this section, we perform a noise sensitivity evaluation in the Franka Kitchen environment to assess the robustness of different methods. Specifically, we adjust the reversal rate r from 50% (as used in previous experiments) to 25% and 75%, to evaluate the method's stability under different levels of inconsistency. For clarity, we present only the most challenging p4 metric.

*Table 3.* Evaluation results of p4 metric under different levels of reverse rates in the Kitchen environment.

|             | r=25%                            | r=50%                            | r=75%                            |
|-------------|----------------------------------|----------------------------------|----------------------------------|
| BET-CPL     | $65.7 \pm 1.6$                   | $62.6\pm2.0$                     | $55.0\pm2.5$                     |
| BET-DPO     | $60.2\pm4.8$                     | $57.4\pm6.6$                     | $47.2\pm7.0$                     |
| Diff-CPL    | $66.0\pm1.0$                     | $63.5\pm0.8$                     | $57.1\pm2.5$                     |
| FKPD        | $71.3\pm2.3$                     | $64.1\pm3.2$                     | $62.3\pm4.6$                     |
| Diff-UAPA-C | $75.3\pm2.9$                     | $70.9\pm2.5$                     | $\textbf{70.5} \pm \textbf{3.8}$ |
| Diff-UAPA-I | $\textbf{75.5} \pm \textbf{3.0}$ | $\textbf{71.7} \pm \textbf{4.6}$ | $69.1\pm5.2$                     |
|             |                                  |                                  |                                  |

**Results Analysis.** Table 3 presents the evaluation results. As the noise level increases (i.e., the reversal rate), all methods show a decline in performance, highlighting the significance of uncertainties in the dataset. However, compared to the other methods, Diff-UAPA consistently exhibits better performance with the highest success rate regardless of the scale of noise. This underscores the effectiveness of incorporating the Beta prior model to handle such uncertainties.

# 6. Limitation

**Offline Trajectory Dataset.** This paper primarily focuses on learning from an offline trajectory dataset with potentially inconsistent human preferences that are iteratively updated, where the agent cannot directly interact with the environment. This partial offline setup may limit the agent's ability to explore and discover improved strategies through interactive online learning. However, our method can also generalize to an online setting, where both trajectories and human preferences are dynamically updated over time.

**Computational Overhead.** The integration of training a Beta prior model through variational inference adds computational complexity compared to simpler MLE-based methods. However, by utilizing efficient techniques like the reparameterization trick to enhance scalability, the computational overhead of training the Beta model is minimal in practice, adding only a small additional time cost relative to the diffusion training process.

#### 7. Conclusion

In this paper, we present an uncertainty-aware preference alignment approach for diffusion policies using an iteratively updated preference dataset. Building on the maximum likelihood objective for directly aligning diffusion policies without learning a reward model, we introduce a Maximum A Posteriori (MAP) objective with an informative Beta prior, which is capable of capturing the uncertainty arising from potentially inconsistent human preferences. Empirical results across various domains demonstrate the effectiveness of our method. For future work, we extend this framework to the online RL setting with complex tasks involving humanoid robots or dexterous hand. By enabling agents to interact with the environment, our system can dynamically adapt to evolving human preferences, thereby solving more difficult applications.

427 428

429

430

385

## 440 Impact Statement

441 The potential broader impact of this work is significant, as 442 it advances the field of human-aligned decision-making in 443 artificial intelligence (AI) and robotics. From an ethical 444 perspective, this work emphasizes reducing bias and incon-445 sistency in preference-based reinforcement learning, which 446 aligns with the principles of fairness and equity in AI. How-447 ever, challenges remain in ensuring the informed collection 448 of preference data and safeguarding against misuse, such 449 as exploiting preference alignment for manipulative or un-450 ethical purposes. Transparency in how user preferences are 451 modeled and incorporated into decision-making policies is 452 crucial to building trust and accountability. 453

454 Future societal consequences may include the development 455 of AI systems that better reflect the diverse needs of global 456 populations, contributing to more personalized and human-457 centric technologies. However, there is also the risk of over-458 reliance on population-level preferences that might inad-459 vertently marginalize minority views or lead to unintended 460 consequences if preferences are improperly interpreted or 461 misaligned with ethical considerations. Addressing these 462 risks requires careful oversight, interdisciplinary collabora-463 tion, and ongoing dialogue with diverse stakeholders. 464

## References

465

466

471

- 467 Ajay, A., Du, Y., Gupta, A., Tenenbaum, J., Jaakkola, T., and
  468 Agrawal, P. Is conditional generative modeling all you
  469 need for decision-making? In *International Conference*470 on *Learning Representations*, 2023.
- 472 An, G., Lee, J., Zuo, X., Kosaka, N., Kim, K.-M., and Song,
  473 H. O. Direct preference-based policy optimization with474 out reward modeling. *Advances in Neural Information*475 *Processing Systems*, 36:70247–70266, 2023.
- Bradley, R. A. and Terry, M. E. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Casper, S., Davies, X., Shi, C., Gilbert, T. K., Scheurer, J., Rando, J., Freedman, R., Korbak, T., Lindner, D., Freire, P., et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- Chen, H., Lu, C., Ying, C., Su, H., and Zhu, J. Offline reinforcement learning via high-fidelity generative behavior modeling. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Chen, L., Bahl, S., and Pathak, D. Playfusion: Skill acquisition via diffusion from language-annotated play. In *Conference on Robot Learning*, pp. 2012–2029, 2023b.

- Chen, Y., Li, H., and Zhao, D. Boosting continuous control with consistency policy. In Autonomous Agents and Multiagent Systems, pp. 335–344, 2024.
- Chi, C., Xu, Z., Feng, S., Cousineau, E., Du, Y., Burchfiel, B., Tedrake, R., and Song, S. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 2023.
- Choi, H., Jung, S., Ahn, H., and Moon, T. Listwise reward estimation for offline preference-based reinforcement learning. In *International Conference on Machine Learning*, 2024.
- Christiano, P. F., Leike, J., Brown, T. B., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, pp. 4299–4307, 2017.
- Dong, Z., Yuan, Y., HAO, J., Ni, F., Mu, Y., ZHENG, Y., Hu, Y., Lv, T., Fan, C., and Hu, Z. Aligndiff: Aligning diverse human preferences via behavior-customisable diffusion model. In *International Conference on Learning Representations*, 2024.
- Fang, B., Jia, S., Guo, D., Xu, M., Wen, S., and Sun, F. Survey of imitation learning for robotic manipulation. *Int. J. Intell. Robotics Appl.*, 3(4):362–369, 2019.
- Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. D4rl: Datasets for deep data-driven reinforcement learning. arXiv preprint arXiv:2004.07219, 2020.
- Gupta, A., Kumar, V., Lynch, C., Levine, S., and Hausman, K. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. *arXiv preprint arXiv:1910.11956*, 2019.
- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pp. 1352–1361, 2017.
- Hejna, J. and Sadigh, D. Inverse preference learning: Preference-based rl without a reward function. *Advances in Neural Information Processing Systems*, 36, 2024.
- Hejna, J., Rafailov, R., Sikchi, H., Finn, C., Niekum, S., Knox, W. B., and Sadigh, D. Contrastive preference learning: Learning from human feedback without rl. In *International Conference on Learning Representations*, 2024.
- Hejna III, D. J. and Sadigh, D. Few-shot preference learning for human-in-the-loop rl. In *Conference on Robot Learning*, pp. 2014–2025, 2023.

- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in Neural Information Process- ing systems*, 33:6840–6851, 2020.
- Hwang, M., Lee, G., Kee, H., Kim, C. W., Lee, K., and Oh,
  S. Sequential preference ranking for efficient reinforcement learning from human feedback. *Advances in Neural Information Processing Systems*, 36:49088–49099, 2023.
- 503 Ibarz, B., Leike, J., Pohlen, T., Irving, G., Legg, S., and
  504 Amodei, D. Reward learning from human preferences
  505 and demonstrations in atari. In *Advances in Neural In-*506 *formation Processing Systems, NeurIPS*, pp. 8022–8034,
  507 2018.
- Janner, M., Du, Y., Tenenbaum, J. B., and Levine, S. Planning with diffusion for flexible behavior synthesis. In *International Conference on Machine Learning*, 2022.
  - Joo, W., Lee, W., Park, S., and Moon, I. Dirichlet variational autoencoder. *Pattern Recognition*, 107:107514, 2020.

- Kang, B., Ma, X., Du, C., Pang, T., and Yan, S. Efficient diffusion policies for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2023a.
- Kang, Y., Shi, D., Liu, J., He, L., and Wang, D. Beyond reward: Offline preference-guided policy optimization. In *International Conference on Machine Learning*, 2023b.
- Kim, C., Park, J., Shin, J., Lee, H., Abbeel, P., and Lee,
  K. Preference transformer: Modeling human preferences
  using transformers for RL. In *International Conference on Learning Representations*, 2023.
- Knox, W. B., Hatgis-Kessell, S., Booth, S., Niekum, S.,
  Stone, P., and Allievi, A. Models of human preference for learning reward functions. *arXiv preprint arXiv:2206.02231*, 2022.
- Lee, K., Smith, L., and Abbeel, P. Pebble: Feedbackefficient interactive reinforcement learning via relabeling
  experience and unsupervised pre-training. In *Interna- tional Conference on Machine Learning*, 2021.
- Liang, X., Shu, K., Lee, K., and Abbeel, P. Reward uncertainty for exploration in preference-based reinforcement learning. *arXiv preprint arXiv:2205.12401*, 2022.
- Liu, R., Bai, F., Du, Y., and Yang, Y. Meta-reward-net: Implicitly differentiable reward learning for preferencebased reinforcement learning. *Advances in Neural Information Processing Systems*, 35:22270–22284, 2022.
- Lu, C., Chen, H., Chen, J., Su, H., Li, C., and Zhu, J. Contrastive energy prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In *International Conference on Machine Learning*, pp. 22825– 22855, 2023.

- MacGlashan, J., Ho, M. K., Loftin, R., Peng, B., Wang, G., Roberts, D. L., Taylor, M. E., and Littman, M. L. Interactive learning from policy-dependent human feedback. In *International Conference on Machine Learning*, pp. 2285–2294, 2017.
- Mandlekar, A., Xu, D., Wong, J., Nasiriany, S., Wang, C., Kulkarni, R., Fei-Fei, L., Savarese, S., Zhu, Y., and Martín-Martín, R. What matters in learning from offline human demonstrations for robot manipulation. In *Conference on Robot Learning (CoRL)*, 2021.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.
- Mukhoti, J., Kulharia, V., Sanyal, A., Golodetz, S., Torr, P., and Dokania, P. Calibrating deep neural networks using focal loss. *Advances in Neural Information Processing Systems*, 33:15288–15299, 2020.
- Newman, M. E. Efficient computation of rankings from pairwise comparisons. *Journal of Machine Learning Research*, 24(238):1–25, 2023.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Park, J., Seo, Y., Shin, J., Lee, H., Abbeel, P., and Lee, K. Surf: Semi-supervised reward learning with data augmentation for feedback-efficient preference-based reinforcement learning. In *International Conference on Learning Representations*, 2022.
- Psenka, M., Escontrela, A., Abbeel, P., and Ma, Y. Learning a diffusion model policy from rewards via q-score matching. In *International Conference on Machine Learning*, 2024.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. In Advances in Neural Information Processing Systems, 2023.
- Rafailov, R., Hejna, J., Park, R., and Finn, C. From \$r\$ to \$q^\*\$: Your language model is secretly a q-function. In *First Conference on Language Modeling*, 2024.
- Shafiullah, N. M. M., Cui, Z. J., Altanzaya, A., and Pinto, L. Behavior transformers: Cloning \$k\$ modes with one stone. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), Advances in Neural Information Processing Systems, 2022.

- Shan, Z., Fan, C., Qiu, S., Shi, J., and Bai, C. Forward kl regularized preference optimization for aligning diffusion policies. *arXiv preprint arXiv:2409.05622*, 2024.
- Shin, D., Dragan, A. D., and Brown, D. S. Benchmarks and algorithms for offline preference-based reward learning. *Transactions on Machine Learning Research*, 2023.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L.,
  Van Den Driessche, G., Schrittwieser, J., Antonoglou, I.,
  Panneershelvam, V., Lanctot, M., et al. Mastering the
  game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Sridhar, A., Shah, D., Glossop, C., and Levine, S. Nomad:
  Goal masked diffusion policies for navigation and exploration. In *IEEE International Conference on Robotics* and Automation, pp. 63–70, 2024.
  - Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction.* MIT press, 2018.

573

574

575

576

577

578

- Vaswani, A. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Verma, M. and Metcalf, K. Hindsight priors for reward learning from human preferences. In *International Conference on Learning Representations*, 2024.
- Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A.,
  Purushwalkam, S., Ermon, S., Xiong, C., Joty, S., and
  Naik, N. Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8228–8238, 2024.
- Wang, X., Wang, S., Liang, X., Zhao, D., Huang, J., Xu,
  X., Dai, B., and Miao, Q. Deep reinforcement learning:
  A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4):5064–5078, 2022.
- Wang, Z., Hunt, J. J., and Zhou, M. Diffusion policies as an expressive policy class for offline reinforcement learning. In *International Conference on Learning Representations*, 2023.
- Warnell, G., Waytowich, N., Lawhern, V., and Stone,
  P. Deep tamer: Interactive agent shaping in highdimensional state spaces. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- Wirth, C., Akrour, R., Neumann, G., and Fürnkranz, J. A
  survey of preference-based reinforcement learning methods. *Journal of Machine Learning Research*, 18(136): 1–46, 2017.

- Xu, M., Xu, Z., Chi, C., Veloso, M., and Song, S. Xskill: Cross embodiment skill discovery. In *Conference on Robot Learning*, pp. 3536–3555, 2023.
- Xu, S., Yue, B., Zha, H., and Liu, G. A distributional approach to uncertainty-aware preference alignment using offline demonstrations. In *International Conference on Learning Representations*, 2025.
- Xue, W., An, B., Yan, S., and Xu, Z. Reinforcement learning from diverse human preferences. In *International Joint Conference on Artificial Intelligence*, 2024.
- Yang, L., Huang, Z., Lei, F., Zhong, Y., Yang, Y., Fang, C., Wen, S., Zhou, B., and Lin, Z. Policy representation via diffusion probability model for reinforcement learning. *arXiv preprint arXiv:2305.13122*, 2023.
- Yuan, Y., Hao, J., Ma, Y., Dong, Z., Liang, H., Liu, J., Feng, Z., Zhao, K., and Zheng, Y. Uni-rlhf: Universal platform and benchmark suite for reinforcement learning with diverse human feedback. In *International Conference on Learning Representations, ICLR*, 2024.
- Zhu, Z., Zhao, H., He, H., Zhong, Y., Zhang, S., Guo, H., Chen, T., and Zhang, W. Diffusion models for reinforcement learning: A survey. arXiv preprint arXiv:2311.01223, 2023.
- Ziebart, B. D. Modeling purposeful adaptive behavior with the principle of maximum causal entropy, 2010.

# A. More Details in Section 4.1

We detailed the deviation from Equation (11) to Equation (12) here.

$$\begin{split} \mathcal{L}_{1,\text{MLE}}^{(\tau^{w},\tau^{l})}(\theta) \\ &= -\log \sigma \Big( \alpha \cdot \Big( \sum_{t=0}^{k} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,w}|s_{t}^{w},a_{t}^{0,w})} \left[ \gamma^{t} \log \frac{\pi_{\theta}(\overline{a_{t}^{w}}|s_{t}^{w})}{\pi_{\text{ref}}(a_{t}^{w}|s_{t}^{w})} \right] - \sum_{t=0}^{k} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,l}|s_{t}^{l},a_{t}^{0,l})} \left[ \gamma^{t} \log \frac{\pi_{\theta}(\overline{a_{t}^{l}}|s_{t}^{l})}{\pi_{\text{ref}}(\overline{a_{t}^{l}}|s_{t}^{l})} \right] \Big) \Big) \\ &= -\log \sigma \Big( \alpha \cdot \Big( \sum_{t=0}^{k} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,r}|s_{t}^{l},a_{t}^{0,r})} \left[ \gamma^{t} \log \frac{\pi_{\theta}(\overline{a_{t}^{w}}|s_{t}^{w})}{\pi_{\text{ref}}(\overline{a_{t}^{w}}|s_{t}^{w})} - \gamma^{t} \log \frac{\pi_{\theta}(\overline{a_{t}^{l}}|s_{t}^{l})}{\pi_{\text{ref}}(\overline{a_{t}^{l}}|s_{t}^{l})} \right] \Big) \Big) \\ &= -\log \sigma \Big( \alpha \cdot \Big( \sum_{t=0}^{k} \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,r}|s_{t}^{l},a_{t}^{0,r})} \left[ \sum_{i=1}^{I} \left( \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w})}{\pi_{\text{ref}}(a_{t}^{i-1|i,w}|s_{t}^{w})} - \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l})}{\pi_{\text{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})} \right) \Big] \Big) \Big) \\ &= -\log \sigma \Big( \alpha \cdot \Big( \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,r}|s_{t}^{l},a_{t}^{0,r})} \left[ \sum_{t=0}^{k} \sum_{i=1}^{I} \left( \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w})}{\pi_{\text{ref}}(a_{t}^{i-1|i,w}|s_{t}^{w})} - \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l})}{\pi_{\text{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})} \right) \Big] \Big) \Big) \\ &= -\log \sigma \Big( \alpha I \cdot \Big( \mathbb{E}_{\pi_{\theta}(a_{t}^{1:I,r}|s_{t}^{l},a_{t}^{0,r})} \pi_{\theta}(a_{t}^{0,w}|s_{t}^{w},a_{t}^{i,w})}, \left[ \sum_{t=0}^{k} \left( \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w})}{\pi_{\text{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})} - \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l})}{\pi_{\text{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})} - \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l})}{\pi_{\text{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})} \right) \Big] \Big) \Big) \\ &= -\log \sigma \Big( \alpha I \cdot \Big( \mathbb{E}_{a_{t}^{i,w} \sim q(a_{t}^{i}|s_{t}^{w},a_{0}^{0,u})\pi_{\theta}(a_{t}^{i-1,l}|s_{t}^{w},a_{t}^{i,v})}, \left[ \sum_{t=0}^{k} \left( \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w})}{\pi_{\text{ref}}(a_{t}^{i-1|i,w}|s_{t}^{w})} - \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l})}{\pi_{\text{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})} \right) \Big] \Big) \Big) \\ &= -\log \sigma \Big( \alpha I \cdot \Big( \mathbb{E}_{a_{t}^{i,w} \sim q(a_{t}^{i}|s_{t}^{w},a_{0}^{0,u})\pi_{\theta}(a_{t}^{i-1,l}|s_{t}^{w},a_{t}^{i,v})}, \left[ \sum_{t=0}^{k} \left( \gamma^{t} \log \frac{\pi_{\theta}(a_{t}^{i-1$$

Since  $-\log \sigma(x)$  is a convex function:

$$(-\log \sigma(x))'' = (\sigma(x) - 1)' = (\sigma(x)(1 - \sigma(x))) \ge 0$$

According to Jensen's inequality:

$$\mathbb{E}_{a_{t}^{i,w} \sim q(a_{t}^{i,w}|a_{t}^{0,w},s_{t}^{w}), a_{t}^{i,l} \sim q(a_{t}^{i,l}|a_{t}^{0,w},s_{t}^{w}), \left[ -\log\sigma\left(\alpha I \cdot \left(\mathbb{E}_{a_{t}^{i-1,\cdot} \sim \pi_{\theta}(a_{t}^{i-1,\cdot}|s_{t}^{\cdot},a_{t}^{0,\cdot})}[\sum_{t=0}^{k}(\gamma^{t}\log\frac{\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w})}{\pi_{\mathrm{ref}}(a_{t}^{i-1|i,w}|s_{t}^{w})} - \gamma^{t}\log\frac{\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l})}{\pi_{\mathrm{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})})]\right)\right) \right] \\ = \mathbb{E}_{a_{t}^{i,w} \sim q(a_{t}^{i,u}|a_{t}^{0,w},s_{t}^{w}), a_{t}^{i,l} \sim q(a_{t}^{i,l}|a_{t}^{0,u},s_{t}^{w}), \left[ -\log\sigma\left(\alpha I \cdot \sum_{t=0}^{k}\left(\gamma^{t}\mathbb{D}_{\mathrm{KL}}\left[\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w}) \mid | \pi_{\mathrm{ref}}(a_{t}^{i-1|i,w}|s_{t}^{w})\right] - \gamma^{t}\mathbb{D}_{\mathrm{KL}}\left[\pi_{\theta}(a_{t}^{i-1|i,l}|s_{t}^{l}) \mid | \pi_{\mathrm{ref}}(a_{t}^{i-1|i,l}|s_{t}^{l})\right] \right)\right) \right]$$

According to Formula (1), it can be further simplified as:

$$-\mathbb{E}_{a_{t}^{i,w} \sim q(a_{t}^{i,w}|a_{t}^{0,w},s_{t}^{w}),} \Big[ \log \sigma \Big( -\alpha I \cdot \Big( \sum_{t=0}^{k} \gamma^{t} (\|\epsilon^{w} - \epsilon_{\theta}(a_{t}^{i,w},s_{t}^{w},i)\|_{2}^{2} - \|\epsilon^{w} - \epsilon_{\text{ref}}(a_{t}^{n,w},s_{t}^{w},i)\|_{2}^{2} \Big) \\ - \sum_{t=0}^{k} \gamma^{t} (\|\epsilon^{l} - \epsilon_{\theta}(a_{t}^{i,l},s_{t}^{l},i)\|_{2}^{2} - \|\epsilon^{l} - \epsilon_{\text{ref}}(a_{t}^{i,l},s_{t}^{l},i)\|_{2}^{2} \Big) \Big) \Big]$$

where 1)  $i \sim \mathcal{U}(0, I)$  is the diffusion timestep, 2)  $a_t^{i,w/l} \sim q(a_t^{i,w/l} | a_t^{0,w/l}, s^{w/l}$  denotes the action  $a_t^{0,w/l}$  corrupted with noise  $\epsilon^{w/l}$  after *i* diffusion steps, and 3)  $\epsilon_{\theta}^{w/l}$  is the noise predictor.

#### **B.** More Details in Section 4.2

We detailed the deviation of Equation (16) here.

and hence

7(

the inner part Beta( $x; \alpha + 1, \beta + 1$ ) serves as the real argument t, and the composition preserves convexity, implying f(x)is convex.

- According to Jensen's inequality
- $\begin{array}{c} 709\\ 710\\ 8a_{t}^{\tau\in\mathcal{D}_{\text{pref}}},\\ 711\\ 712 \end{array} \left[ -\log\sigma \Big( \text{Beta}\Big(\alpha I \cdot \Big( \mathbb{E}_{a_{t}^{i-1,\cdot} \sim \pi_{\theta}(a_{t}^{i-1,\cdot}|s_{t}^{*},a_{t}^{0},\cdot)} [(\sum_{t=0}^{k}\gamma^{t}\log\frac{\pi_{\theta}(a_{t}^{i-1|i,w}|s_{t}^{w})}{\pi_{\text{ref}}(a_{t}^{i-1|i,w}|s_{t}^{w})} \sum_{\tau\in\mathcal{D}_{\text{pref}},t=0}^{k}\frac{\gamma^{t}}{|\mathcal{D}_{\text{pref}}|}\log\frac{\pi_{\theta}(a_{t}^{i-1|i,\tau}|s_{t}^{\tau})}{\pi_{\text{ref}}(a_{t}^{i-1|i,\tau}|s_{t}^{\tau})})]\big);\alpha+1,\beta+1\Big)\Big) \right]$

$$\frac{1}{713} = \mathbb{E}_{\substack{\tau \in \mathcal{D}_{\text{pref}}, \\ \tau \neq a_t^{i, \cdot} \sim q(a_t^{i, \cdot} \mid a_t^{0, \cdot}, s_t^{i})}} \left[ -\log \sigma \left( \alpha I \cdot \sum_{t=0}^k \left( \gamma^t \mathbb{D}_{\text{KL}} \left[ \pi_{\theta}(a_t^{i-1 \mid i, w} \mid s_t^w) \mid \mid \pi_{\text{ref}}(a_t^{i-1 \mid i, w} \mid s_t^w) \right] - \sum_{\tau \in \mathcal{D}_{\text{pref}}}^k \frac{\gamma^t}{\mid \mathcal{D}_{\text{pref}} \mid} \mathbb{D}_{\text{KL}} \left[ \pi_{\theta}(a_t^{i-1 \mid i, \tau} \mid s_t^\tau) \mid \mid \pi_{\text{ref}}(a_t^{i-1 \mid i, \tau} \mid s_t^\tau) \right] \right) \right) \right]$$

$$\begin{aligned} \mathcal{L}_{1,\text{prior}}^{(\tau^w,\tau^l)}(\theta) \\ &= -\log\sigma\Big(\text{Beta}\Big(\alpha \cdot \Big(\sum_{t=0}^k \mathbb{E}_{\pi_{\theta}(a_t^{1:I,w}|s_t^w,a_t^{0,w})} \left[\gamma^t \log \frac{\pi_{\theta}(\overline{a_t^w}|s_t^w)}{\pi_{\text{ref}}(a_t^w|s_t^w)}\right] \\ &- \sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \mathbb{E}_{\pi_{\theta}(a_t^{1:I,\tau}|s_t^{,\tau},a_t^{0,\tau})} \left[\gamma^t \log \frac{\pi_{\theta}(\overline{a_t^w}|s_t^v)}{\pi_{\text{ref}}(\overline{a_t^w}|s_t^w)}\right]; \alpha + 1, \beta + 1\Big)\Big) \\ &= -\log\sigma\Big(\text{Beta}\Big(\alpha \cdot \Big(\mathbb{E}_{\pi_{\theta}(a_t^{1:I,\tau}|s_t,a_t^{0,\tau})} \left[\sum_{t=0}^k \gamma^t \log \frac{\pi_{\theta}(\overline{a_t^w}|s_t^w)}{\pi_{\text{ref}}(\overline{a_t^w}|s_t^w)} - \sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \frac{\gamma^t}{|\mathcal{D}_{\text{pref}}|} \log \frac{\pi_{\theta}(\overline{a_t^w}|s_t^\tau)}{\pi_{\text{ref}}(\overline{a_t^v}|s_t^\tau)}\right]; \alpha + 1, \beta + 1\Big)\Big) \\ &= -\log\sigma\Big(\alpha \cdot \Big(\sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \mathbb{E}_{\pi_{\theta}(a_t^{1:I,\tau}|s_t,a_t^{0,\tau})} \left[\sum_{i=1}^I \left(\sum_{t=0}^k \gamma^t \log \frac{\pi_{\theta}(a_t^{i-1|i,w}|s_t^w)}{\pi_{\text{ref}}(a_t^{i-1|i,w}|s_t^w)} - \sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \frac{\gamma^t}{|\mathcal{D}_{\text{pref}}|} \log \frac{\pi_{\theta}(a_t^{i-1|i,\tau}|s_t^\tau,s_t^\tau)}{\pi_{\text{ref}}(a_t^{i-1|i,\tau}|s_t^\tau)}\Big)\Big]\Big)\Big) \\ &= -\log\sigma\Big(\alpha \cdot \Big(\mathbb{E}_{\pi_{\theta}(a_t^{::I,\tau}|s_t,a_t^{0,\tau})} \left[\left(\sum_{t=0}^k \sum_{i=1}^I \gamma^t \log \frac{\pi_{\theta}(a_t^{i-1|i,w}|s_t^w)}{\pi_{\text{ref}}(a_t^{i-1|i,w}|s_t^w)} - \sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \sum_{i=1}^I \frac{\gamma^t}{|\mathcal{D}_{\text{pref}}|} \log \frac{\pi_{\theta}(a_t^{i-1|i,\tau}|s_t^\tau)}{\pi_{\text{ref}}(a_t^{i-1|i,\tau}|s_t^\tau)}\Big)\Big]\Big)\Big) \\ &= -\log\sigma\Big(\alpha I \cdot \Big(\mathbb{E}_{a_t^{i,\cdots},\alpha(a_t^{i}|s_t,a_t^{0,\cdots})} \left[\left(\sum_{t=0}^k \gamma^t \log \frac{\pi_{\theta}(a_t^{i-1|i,w}|s_t^w)}{\pi_{\text{ref}}(a_t^{i-1|i,w}|s_t^w)} - \sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \frac{\gamma^t}{|\mathcal{D}_{\text{pref}}|} \log \frac{\pi_{\theta}(a_t^{i-1|i,\tau}|s_t^\tau)}{\pi_{\text{ref}}(a_t^{i-1|i,\tau}|s_t^\tau)}\Big)\right]\Big)\Big) \\ &= -\log\sigma\Big(\alpha I \cdot \Big(\mathbb{E}_{a_t^{i,\cdots},\alpha(a_t^{i}|s_t,a_t^{0,\cdots})}\pi_{\theta}(a_t^{0,\cdots}|s_t,a_t^{i,\cdots})} \left[\left(\sum_{t=0}^k \gamma^t \log \frac{\pi_{\theta}(a_t^{i-1|i,w}|s_t^w)}{\pi_{\text{ref}}(a_t^{i-1|i,w}|s_t^w)} - \sum_{\tau \in \mathcal{D}_{\text{pref},t=0}}^k \frac{\gamma^t}{|\mathcal{D}_{\text{pref}}|}\log \frac{\pi_{\theta}(a_t^{i-1|i,\tau}|s_t^\tau)}{\pi_{\text{ref}}(a_t^{i-1|i,\tau}|s_t^\tau)}\Big)\Big]\Big)\Big) \end{aligned}$$

Since  $-\log \sigma(\text{Beta}(x; \alpha, \beta))$  is a convex function when  $\alpha + \beta \ge 2$ . Define  $g(t) = -\log(\sigma(t))$ . Since

 $-\log(\sigma(t)) = \log(1 + e^{-t}),$ 

it suffices to show that  $\log(1 + e^{-t})$  is convex in t. Differentiating,

This shows  $\log(1 + e^{-t})$  is strictly convex in t. Therefore, for the function

$$\frac{d}{dt}\log(1+e^{-t}) = \frac{-e^{-t}}{1+e^{-t}} = -\frac{1}{e^t+1}$$

 $\frac{d^2}{dt^2} \log(1 + e^{-t}) = \frac{e^t}{(e^t + 1)^2} > 0 \quad (\forall t \in \mathbb{R}).$ 

 $f(x) = -\log \left[ \sigma \left( \operatorname{Beta}(x; \alpha + 1, \beta + 1) \right) \right],$ 

According to Formula (1), it can be further simplified as:

$$-\mathbb{E}_{\substack{a_t^{i,\cdot} \sim q(a_t^{i,\cdot} \mid a_t^{0,\cdot}, s_t^{\cdot})}} \left[ \log \sigma \Big( -\alpha I \cdot \Big( \sum_{t=0}^k \gamma^t (\|\epsilon^w - \epsilon_\theta(a_t^{i,w}, s_t^w, i)\|_2^2 - \|\epsilon^w - \epsilon_{\text{ref}}(a_t^{n,w}, s_t^w, i)\|_2^2 \right) \\ - \sum_{\tau \in \mathcal{D}_{\text{pref}}, t=0}^k \frac{\gamma^t}{|\mathcal{D}_{\text{pref}}|} (\|\epsilon^\tau - \epsilon_\theta(a_t^{i,\tau}, s_t^\tau, i)\|_2^2 - \|\epsilon^\tau - \epsilon_{\text{ref}}(a_t^{i,\tau}, s_t^\tau, i)\|^2) \Big) \Big) \Big]$$

where 1)  $i \sim \mathcal{U}(0, I)$  is the diffusion timestep, 2)  $a_t^{i,\cdot} \sim q(a_t^{i,\cdot} | a_t^{0,\cdot}, s^{\cdot})$  denotes the action  $a_t^{0,\cdot}$  corrupted with noise  $\epsilon^{\cdot}$  after *i* diffusion steps, and 3)  $\epsilon_{\theta}^{\cdot}$  is the noise predictor.

#### C. Proof of Proposition 4.1

Proposition 4.1 can be divided into two parts: 1) the uncertainty-aware property of the Beta prior, and 2) the prior on the strength of a trajectory.

**Part 1.** We show the uncertainty-aware capability of the Beta prior  $Beta(\phi(\tau); \alpha, \beta)$  during the iterative preference alignment process outlined in Definition 3.1 as follows.

The probability density function (PDF) of the Beta distribution  $\text{Beta}(\phi(\tau); \alpha, \beta)$  is given by:

$$f(\phi(\tau); \alpha, \beta) = \frac{\phi(\tau)^{\alpha - 1} (1 - \phi(\tau))^{\beta - 1}}{B(\alpha, \beta)}, \quad 0 \le \phi(\tau) \le 1,$$
(19)

where  $B(\alpha, \beta) = \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt$  is the Beta function, serving as a normalizing constant.

The variance of a Beta distribution  $\text{Beta}(\phi(\tau); \alpha, \beta)$  is given by the following formula:

$$\operatorname{Var}(\operatorname{Beta}(\alpha,\beta)) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}.$$
(20)

In the process described in Definition 3.1, the uncertainty arises from the varying preferences of different human raters for a given trajectory pair  $(\tau^i, \tau^j)$ . Without loss of generality, assuming an initial belief of Beta(1,1) for each trajectory, and with 10 raters evaluating a candidate pair  $(\tau^i, \tau^j)$ , the Beta prior is updated according to the preferences expressed by the raters. For instance, in the first case, where 9 raters prefer  $\tau^i$  and 1 rater prefers  $\tau^j$ , the Beta prior for  $\tau^i$  would be updated to Beta(10, 2). In the second case, where 5 raters prefer  $\tau^i$  and 5 prefer  $\tau^j$ , the Beta prior for  $\tau^i$  would be Beta(6, 6). Intuitively, we would be more confident with less uncertainty in the first case, as the majority of raters share the same preference.

The Beta distribution effectively captures this uncertainty. As shown in Equation (20), the variance of Beta(10, 2) is smaller than that Beta(6, 6), indicating that Beta(10, 2) is 'sharper' and reflects less uncertainty, which aligns with our intuition.

Part 2. We prove that the prior on the strength of a trajectory is proportional to  $\text{Beta}((\phi(\tau); \alpha+1, \beta+1))$ , i.e.,  $p_0(A^{\pi_{\theta}}(\tau)) \propto \text{Beta}(\phi(\tau); \alpha+1, \beta+1)$ , as follows.

Recall that the probability of a trajectory  $\tau$  with strength  $A^{\pi_{\theta}}(\tau)$  winning against the average candidate is given by  $\phi(\tau) = \sigma(A^{\pi_{\theta}}(\tau) - \bar{A}^{\pi_{\theta}}) \in (0, 1)$ . Let  $A^{\pi_{\theta}}(\tau) - \bar{A}^{\pi_{\theta}}$  be denoted as  $\tilde{A}^{\pi_{\theta}}(\tau)$ . According to Equation (19), we have that the Beta distribution over  $\phi(\tau) = \sigma(\tilde{A}^{\pi_{\theta}}(\tau))$  is:

$$\operatorname{Beta}(\sigma(\tilde{A}^{\pi_{\theta}}(\tau)); \alpha, \beta) \propto \sigma(\tilde{A}^{\pi_{\theta}}(\tau))^{\alpha-1} (1 - \sigma(\tilde{A}^{\pi_{\theta}}(\tau)))^{\beta-1}.$$
(21)

The derivative of the sigmoid function is:

$$\sigma'(\tilde{A}^{\pi_{\theta}}(\tau)) = \sigma(\tilde{A}^{\pi_{\theta}}(\tau))(1 - \sigma(\tilde{A}^{\pi_{\theta}}(\tau))).$$
(22)

770 By incorporating Equation (21) and Equation (22) into Equation (15), we have that:

$$p_{0}(A^{\pi_{\theta}}(\tau)) \propto \sigma(\tilde{A}^{\pi_{\theta}}(\tau))^{\alpha} (1 - \sigma(\tilde{A}^{\pi_{\theta}}(\tau)))^{\beta}$$
  

$$\propto \operatorname{Beta}(\sigma(\tilde{A}^{\pi_{\theta}}(\tau)); \alpha + 1, \beta + 1)$$
  

$$= \operatorname{Beta}(\phi(\tau); \alpha + 1, \beta + 1).$$
(23)

# **D.** More Experimental Details

#### D.1. Experimental Settings

In this paper, we utilized a total of 4 NVIDIA GeForce RTX 3090 GPUs, each with 24 GB of memory. The random seeds used for the experiments were 42, 43, and 44. We trained the agents offline and selected the final epoch for evaluation across 56 parallel environments, each with 10 episodes. Additionally, we employed a transformer-based architecture for the Beta model as in the preference transformer (Kim et al., 2023).

## **D.2.** Environmental Details

**Manipulation Tasks.** Robomimic (Mandlekar et al., 2021) is a large-scale robotic manipulation benchmark designed to explore imitation learning and offline reinforcement learning (RL). It consists of five tasks, each with a proficient human (PH) teleoperated demonstration dataset, and four tasks also feature mixed proficient/non-proficient human (MH) demonstration datasets, resulting in a total of nine variants. In this paper, we focus on three tasks: Lift, Can, and Square. Specifically:

- Lift: The robot arm must lift a small cube. This is the simplest task.
- Can: The robot must move a Coke can from a large bin to a smaller target bin. This task is slightly more challenging than Lift, as picking up the can is more difficult than picking up the cube, and the can must be placed accurately in the target bin.
- Square: The robot is required to pick up a square nut and place it onto a rod. This task is significantly more difficult than Lift and Can, as it demands high precision to pick up the nut and insert it into the rod.

The Franka Kitchen is also a widely used environment for evaluating the performance of methods in learning complex, long-horizon tasks. Introduced in Relay Policy Learning (Gupta et al., 2019), the environment features seven objects for interaction and includes a human demonstration dataset consisting of 566 demonstrations, each completing four tasks in random order. The objective is to execute as many of the demonstrated tasks as possible, regardless of their order, highlighting both short-horizon and long-horizon multimodal capabilities.

**Locomotion Tasks.** We evaluate our locomotion tasks using the D4RL benchmark (Fu et al., 2020), which is widely used in reinforcement learning (RL) for continuous control tasks. In this paper, we focus on the Hopper and HalfCheetah environments. In these environments, the goal is to maximize the cumulative reward within a single episode by navigating a sequence of actions that optimize the agent's movement and efficiency. More specifically:

- Hopper: In this task, the agent controls a 2D hopping robot, with the objective of balancing and moving the robot forward using as few steps as possible.
- HalfCheetah: In this task, the agent controls a 2D robotic cheetah, aiming to run as fast as possible while maximizing speed and maintaining stability.

# D.3. Manipulation Preference Dataset

For the robot manipulation tasks, we train two policies using behavior cloning: the BET policy and the diffusion policy. Training proceeds until a 40% success rate is reached. To build the simulation environment, we deploy 56 parallel environments, each initialized with a different seed to ensure varied initial positions for the agent. We then collect 560 trajectories per policy. From these, we randomly select 500 trajectory pairs and label them based on the sum of their rewards. During training, each trajectory is sliced using the observed steps as the stride, and these segments are compared. In the iterative update process, for each update round, we randomly select 20% of the trajectory pairs and apply a 50% reversal rate by swapping the winner and loser. To improve stability and convergence, the learning rate is reset at the start of each round.

#### **D.4.** Hyperparameters

Our experiments are primarily based on the codebase from (Chi et al., 2023). Therefore, we retain the same hyperparameters for training the diffusion policy as specified in (Chi et al., 2023) for each experiment. The specific hyperparameters for Diff-UAPA are listed in Table 4.

Table 4. List of the specific hyperparameters for the proposed Diff-UAPA. To ensure fair comparisons, we maintain consistency in other parameters of the same neural networks across different models. 

| Parameters                | Robomimic        | Kitchen          | D4RL             |
|---------------------------|------------------|------------------|------------------|
| General                   |                  |                  |                  |
| Training Epochs           | 600              | 600              | 600              |
| Episode Length            | 400              | 280              | 1000             |
| Beta Model                |                  |                  |                  |
| Network                   | 256              | 256              | 256              |
| Learning Rate             | 2e-5             | 2e-5             | 3e-5             |
| Number of Attention Heads | 4                | 4                | 4                |
| Number of Layers          | 2                | 2                | 1                |
| Batch Size                | 32               | 32               | 64               |
| Initial Belief            | $\alpha=\beta=1$ | $\alpha=\beta=1$ | $\alpha=\beta=1$ |
|                           | ,                | ,                | ,                |

## **D.5.** Visualization Results

Figure 2 presents visualization results from the manipulation tasks. It is evident that the baseline method, Diff-CPL, which is trained using the MLE objective, struggles to handle certain critical scenarios, particularly those involving noisy preferences.



Figure 2. Visualization results in four manipulation tasks.