# USING HUMAN FEEDBACK TO FINE-TUNE DIFFUSION MODELS WITHOUT ANY REWARD MODEL

**Kai Yang**[1]* **Jian Tao**[1]* **Jiafei Lyu**[1] **Chunjiang Ge**[2]
**Jiaxin Chen**[3] **Qimai Li**[3] **Weihan Shen**[3] **Xiaolong Zhu**[3] **Xiu Li**[1]
[1]Shenzhen International Graduate School, Tsinghua University
[2]Department of Automation, Tsinghua University
[3]Parametrix Technology Company Ltd.
`{yk22@mails.tsinghua.edu.cn,tj22@mails.tsinghua.edu.cn}`

## ABSTRACT

Using reinforcement learning with human feedback (RLHF) has shown significant promise in fine-tuning diffusion models. Previous methods start by training a reward model that aligns with human preferences, then leverage RL techniques to fine-tune the underlying models. However, crafting an efficient reward model demands extensive datasets, optimal architecture, and manual hyperparameter tuning, making the process both time and cost-intensive. The direct preference optimization (DPO) method, effective in fine-tuning large language models, eliminates the necessity for a reward model. However, the extensive GPU memory requirement of the diffusion model's denoising process hinders the direct application of the DPO method. To address this issue, we introduce the Direct Preference for Denoising Diffusion Policy Optimization (D3PO) method to directly fine-tune diffusion models. The theoretical analysis demonstrates that although D3PO omits training a reward model, it effectively functions as the optimal reward model trained using human feedback data to guide the learning process. This approach requires no training of a reward model, proving to be more direct, cost-effective, and minimizing computational overhead. In experiments, our method uses the relative scale of objectives as a proxy for human preference, delivering comparable results to methods using ground-truth rewards. Moreover, D3PO demonstrates the ability to reduce image distortion rates and generate safer images, overcoming challenges lacking robust reward models. Our code is publicly available in https://github.com/yk7333/D3PO/tree/main

## 1 Introduction

Recent advances in image generation models have yielded unprecedented success in producing high-quality images from textual prompts [1, 2, 3]. Diverse approaches, including Generative Adversarial Networks (GANs) [4], autoregressive models [5, 1, 6, 7, 8, 9], Normalizing Flows [10, 11], and diffusion-based techniques [12, 13, 14, 2], have rapidly pushed forward the capabilities of these systems. With the proper textual inputs, such models are now adept at crafting images that are not only visually compelling but also semantically coherent, garnering widespread interest for their potential applications and implications.

To adapt the aforementioned models for specific downstream tasks, such as the generation of more visually appealing and aesthetic images, Reinforcement Learning from Human Feedback (RLHF) is commonly employed [15]. This technique has been successfully used to refine large language models such as GPT [16, 17]. The method is now being extended to diffusion models to enhance their performance. One such method, the Reward Weighted Regression [18], leverages RLHF to better align generated images with their prompts. An extension of this, the DDPO approach [19], further refines this alignment and seeks to improve image complexity, aesthetic quality, and the congruence between prompt and image. The ReLF technique [20] introduces a novel reward model, dubbed ImageReward, which is specifically trained to discern human aesthetic preferences in text-to-image synthesis. This model is then utilized to fine-tune diffusion models to produce images that align more closely with human preferences. Nonetheless, developing

---

[1]*The first two authors contribute equally to this work.

a robust reward model for various tasks can be challenging, often necessitating a vast collection of images and abundant training resources. For example, to diminish the rate of deformities in character images, one must amass a substantial dataset of deformed and non-deformed images generated from identical prompts. Subsequently, a network is constructed to discern and learn the human preference for non-deformed imagery, serving as the reward model.

In the field of natural language processing, Direct Preference Optimization (DPO) has been proposed to reduce training costs [21]. This method forgoes the training of a reward model and directly fine-tunes language models according to human preferences. However, this straightforward and easy-to-train method encounters challenges when applied to fine-tune diffusion models. During the DPO training process, the complete sentence generated by the language model is treated as a single output, necessitating the storage of gradients from multiple forward passes. With diffusion models, one must store the gradients across multiple latent image representations, which are significantly larger than word embeddings, leading to memory consumption that is typically unsustainable.

To address the issue of high computational overhead and enable the use of the DPO method to fine-tune diffusion models directly with human feedback, we conceptualize the denoising process as a multi-step MDP, which utilizes a pre-trained model to represent an action value function $Q$ that is commonly estimated in RL. We extend the theoretical framework of DPO into the formulated MDP, which allows for direct parameter updates at each step of the denoising process based on human feedback, thereby circumventing the significant computational costs and eliminating the need for a reward model. To the best of our knowledge, this is the first work that forgoes the reward model to fine-tune diffusion models.

Our main contributions are as follows:

- We introduce an innovative approach for fine-tuning diffusion models that could significantly modify the current RLHF framework for fine-tuning diffusion models. This method bypasses resource-intensive reward model training by utilizing direct human feedback, making the process more efficient and cost-effective.
- We expand the theoretical framework of DPO into a multi-step MDP, demonstrating that directly updating the policy based on human preferences within an MDP is equivalent to learning the optimal reward model first and then using it to guide policy updates. This establishes a robust theoretical foundation and provides assurance for our proposed method.
- In our experiments, we have demonstrated the effectiveness of our method by using human feedback to successfully address issues of hand and full-body deformities, enhance the safety of generated images, and improve prompt-image alignment.

## 2  Related Work

**Diffusion models.** Denoising diffusion probabilistic models have quickly become a focal point in the field of generative models, given their remarkable capability to synthesize diverse data types. Originally introduced in [22] and further advanced by [23], these models have been effectively applied to the generation of intricate images [1, 2], dynamic video sequences [24, 25], and complex robotics control systems [26, 27, 28]. The test-to-image diffusion models, in particular, have made it possible to create highly realistic images from textual descriptions [1, 2], opening new possibilities in digital art and design.

In pursuit of more precise manipulation over the generative process, recent studies have explored various techniques to refine the guidance of diffusion models. Adapters, for instance, have been introduced as a means to incorporate additional input constraints, allowing for more targeted generation that adheres to specific criteria [29]. Furthermore, compositional approaches have been developed to blend multiple models for enhanced image quality and generation control [30, 31]. The implementation of classifier [32] and classifier-free guidance [33] has also been a significant step towards achieving more autonomy in the generation process, enabling models to produce outputs that are not only high in fidelity but also align more closely with user intentions. In our work, we use Stable Diffusion [14] to generate images with some specific prompts.

**RLHF.** RLHF stands as a salient strategy in the domain of machine learning when objectives are complex or difficult to define explicitly. This technique has been instrumental across various applications, from gaming, as demonstrated with Atari [15, 34], to more intricate tasks in robotics [35, 36]. The integration of RLHF into the development of large language models (LLMs) has marked a significant milestone in the field, with notable models like OpenAI's GPT-4 [17], Anthropic's Claude [37], Google's Bard [38], and Meta's Llama 2-Chat [39] leveraging this approach to enhance their performance and relevance. The effectiveness of RLHF in refining the behavior of LLMs to be more aligned with human values, such as helpfulness and harmlessness, has been extensively studied [34, 35]. The technique has also proven beneficial in more focused tasks, such as summarization, where models are trained to distill extensive information into concise representations [40].
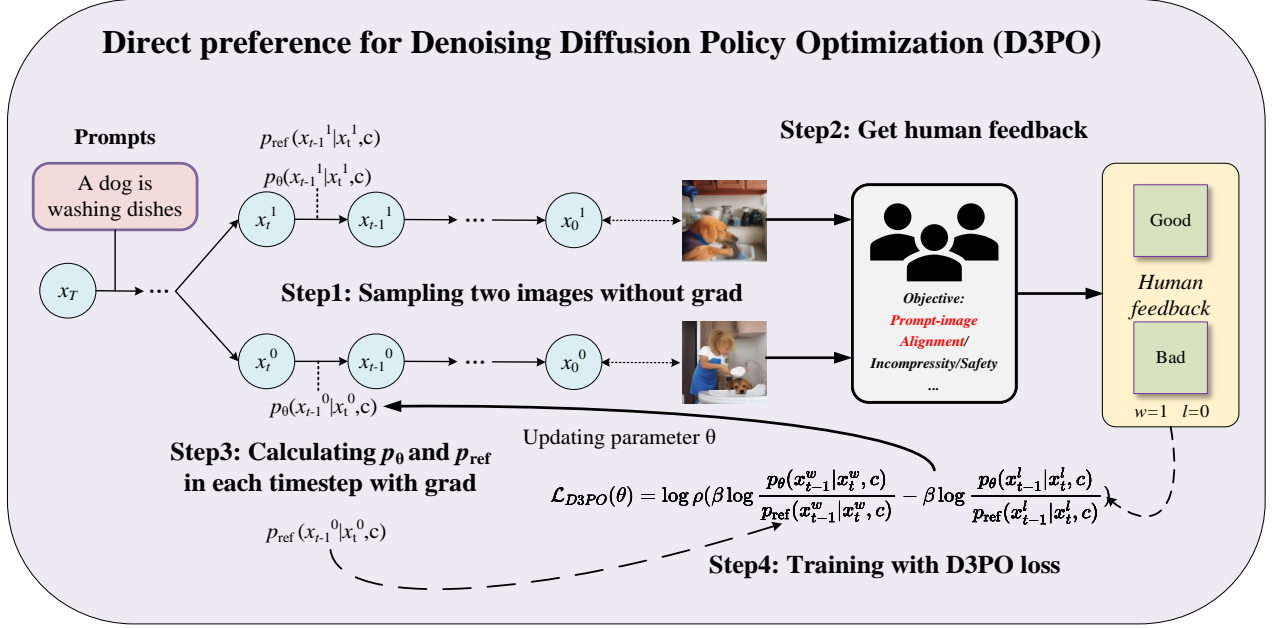
Figure 1: Overview of D3PO. The diffusion model generates two corresponding images based on the provided prompts. Guided by specific task requirements—such as improving prompt-image alignment, enhancing image incompressibility, or refining aesthetic quality—human evaluators select the preferred image. Leveraging this human feedback, our method directly updates the diffusion model's parameters without necessitating the training of a reward model.

**Fine-tune Diffusion Models with RL.** Before applying diffusion models, data generation has been regarded as a sequential decision-making problem and combined with reinforcement learning [41]. More recently, [42] applied reinforcement learning to diffusion models to enhance existing fast DDPM samplers [23]. Reward Weighted method [18] explored using human feedback to align text-to-image models with RLHF. It uses the reward model for the coefficients of the loss function instead of using the reward model when constructing the dataset. ReFL [20] first trains a reward model named ImageReward based on human preferences and then employs it for fine-tuning. During fine-tuning, it randomly selects timesteps to predict the final image, aiming to stabilize the training process by preventing it from solely processing the last denoising step. DDPO [19] treats the denoising process of diffusion models as a MDP to fine-tune diffusion models with many reward models. DPOK [43] combine the KL divergence into the DDPO loss and use it to better align text-to-image objectives. All these models need a robust reward model, which demands a substantial dataset of images and extensive human evaluations.

**Direct Preference Optimization.** In the realm of reinforcement learning, the concept of deriving policies based on preferences, as opposed to explicit rewards, has been explored through various methods. The Contextual Dueling Bandit (CDB) framework [44, 45] replaces the traditional aim for an optimal policy with the concept of a *von Neumann winner*, thereby circumventing the direct use of rewards. *Preference-based Reinforcement Learning* (PbRL) [46, 47, 48] is predicated on learning from binary preferences inferred from an enigmatic *scoring* function instead of explicit rewards. More recently, the work of [21] introduced the DPO approach, which directly utilizes preferences for fine-tuning LLMs. DPO capitalizes on a correlation between the reward function and optimal policies, showing that this constrained reward maximization challenge can be effectively addressed through a singular phase of policy training.

## 3 Preliminaries

**MDP.** We consider the MDP formulation described in [49]. In this setting, an agent perceives a state $s \in \mathcal{S}$ and executes an action, where $\mathcal{S}, \mathcal{A}$ denote state space and action space, respectively. The transition probability function, denoted by $P(s'|s, a)$, governs the progression from the current state $s$ to the subsequent state $s'$ upon the agent's action $a$. Concurrently, the agent is awarded a scalar reward $r$, determined by the reward function $r : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$. The agent's objective is to ascertain a policy $\pi(a|s)$ that maximizes the cumulative returns of trajectories $\tau =$

$(s_0, a_0, s_1, a_1, ..., s_{T-1}, a_{T-1})$, which can be represented as:

$$\mathcal{J}(\pi) = \mathbb{E}_\tau \left[ \sum_{t=0}^{T-1} r\left(s_t, a_t\right) \right]. \tag{1}$$

**Diffusion models.** Diffusion models learn to model a probability distribution $p(x)$ by inverting a Markovian forward process $q(\boldsymbol{x}_t | \boldsymbol{x}_{t-1})$ which adds noise to the data. The denoising process is modeled by a neural network to predict the mean of $\boldsymbol{x}_{t-1}$ or the noise $\epsilon_{t-1}$ of the forward process. In our work, we use network $\boldsymbol{\mu}_\theta(\boldsymbol{x}_t; t)$ to predict the mean of $\boldsymbol{x}_{t-1}$ instead of predicting the noise. Using the mean squared error (MSE) as a measure, the objective of this network can be written as:

$$\mathcal{L} = \mathbb{E}_{t \sim [1,T], \boldsymbol{x}_0 \sim p(\boldsymbol{x}_0), \boldsymbol{x}_t \sim q(\boldsymbol{x}_t | \boldsymbol{x}_0)}[\|\tilde{\boldsymbol{\mu}}(\boldsymbol{x}_0, \boldsymbol{x}_t) - \boldsymbol{\mu}_\theta\left(\boldsymbol{x}_t, t\right)\|^2], \tag{2}$$

where $\tilde{\boldsymbol{\mu}}_\theta(\boldsymbol{x}_t, \boldsymbol{x}_0)$ is the forward process posterior mean. In the case of conditional generative modeling, the diffusion models learn to model $p(x|\boldsymbol{c})$, where $\boldsymbol{c}$ is the conditioning information, i.e., image category and image caption. This is done by adding an additional input $\boldsymbol{c}$ to the denoising neural network, as in $\boldsymbol{\mu}_\theta(\boldsymbol{x}_t, t; \boldsymbol{c})$. To generate a sample from the learned distribution $p_\theta(x|\boldsymbol{c})$, we start by drawing a sample $\boldsymbol{x}_T \sim \mathcal{N}(\boldsymbol{0}, \mathbf{I})$ and then progressively denoise the sample by iterated application of $\epsilon_\theta$ according to a specific sampler [23]. Given the noise-related parameter $\sigma^t$, the reverse process can be written as:

$$p_\theta\left(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{c}\right) = \mathcal{N}\left(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_\theta\left(\boldsymbol{x}_t, \boldsymbol{c}, t\right), \sigma_t^2 \mathbf{I}\right). \tag{3}$$

**Reward learning for preferences.** The basic framework to model preferences is to learn a reward function $r^*(s, a)$ from human feedback [50, 51, 15]. The segment $\sigma = \{s_k, a_k, s_{k+1}, a_{k+1}, ..., s_m, a_m\}$ is a sequence of observations and actions. By following the Bradley-Terry model [52], the human preference distribution $p^*$ by using the reward function can be expressed as:

$$p^*(\sigma_1 \succ \sigma_0) = \frac{\exp(\sum_{t=k}^T r^*(s_t^1, a_t^1))}{\sum_{i \in \{0,1\}} \exp(\sum_{t=k}^T r^*(s_t^i, a_t^i))}, \tag{4}$$

where $\sigma_i \succ \sigma_j$ denotes that segment $i$ is preferable to segment $j$. Now we have the preference distribution of human feedback, and we want to use a network $r_\phi$ to approximate $r^*$. Given the human preference $y \in \{(1,0), (0,1)\}$ which is recorded in dataset $\mathcal{D}$ as a triple $(\sigma^0, \sigma^1, y)$, framing the problem as a binary classification, the reward function modeled as a network is updated by minimizing the following loss:

$$\mathcal{L}(\phi) = -\mathbb{E}_{(\sigma^1, \sigma^0, y) \sim \mathcal{D}}[y(0) \log p_\phi(\sigma_0 \succ \sigma_1) + y(1) \log p_\phi(\sigma_1 \succ \sigma_0)]. \tag{5}$$

## 4 Method

In this section, we describe a method to directly fine-tune diffusion models using human feedback, bypassing the conventional requirement for a reward model. Initially, we reinterpret the denoising process inherent in diffusion models as a multi-step MDP. Then we extend the theory of DPO to MDP, which allows us to apply the principles to effectively translate human preferences into policy improvements in diffusion models.

### 4.1 Denoising process as a multi-step MDP

We conceptualize the denoising process within the diffusion model as a multi-step MDP, which varies slightly from the approach outlined in [19]. To enhance clarity, we have redefined the states, transition probabilities, and policy functions. The correspondence between notations in the diffusion model and the MDP is established as follows:

$$\mathbf{s}_t \triangleq (\boldsymbol{c}, t, \boldsymbol{x}_{T-t}) \qquad\qquad P\left(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t\right) \triangleq \left(\delta_{\boldsymbol{c}}, \delta_{t+1}, \delta_{\boldsymbol{x}_{T-1-t}}\right)$$

$$\mathbf{a}_t \triangleq \boldsymbol{x}_{T-1-t} \qquad\qquad \pi\left(\mathbf{a}_t \mid \mathbf{s}_t\right) \triangleq p_\theta\left(\boldsymbol{x}_{T-1-t} \mid \boldsymbol{x}_{T-t}, \boldsymbol{c}\right)$$

$$\rho_0\left(\mathbf{s}_0\right) \triangleq (p(\boldsymbol{c}), \delta_0, \mathcal{N}(\boldsymbol{0}, \mathbf{I})) \qquad\qquad r(\mathbf{s}_t, \mathbf{a}_t) \triangleq r((\boldsymbol{c}, t, \boldsymbol{x}_{T-t}), \boldsymbol{x}_{T-t-1})$$

where $\delta_x$ represents the Dirac delta distribution, and $T$ denotes the maximize denoising timesteps. Leveraging this mapping, we can employ RL techniques to fine-tune diffusion models by maximizing returns. However, this approach requires a proficient reward model capable of adequately rewarding the noisy images. The task becomes exceptionally challenging, particularly when $t$ is low, and $\boldsymbol{x}_{T-t}$ closely resembles Gaussian noise, even for an experienced expert.

## 4.2 Direct Preference Optimization for MDP

The DPO method does not train a separate reward model but instead directly optimizes the LLMs with the preference data. Given a prompt $x$ and a pair of answers $(y_1, y_0) \sim \pi_{\text{ref}}(y|x)$, where $\pi_{\text{ref}}$ represents the reference (pre-trained) model, these responses are ranked and stored in $\mathcal{D}$ as a tuple $(x, y_w, y_l)$, where $y_w$ denotes the preferred answer and $y_l$ indicates the inferior one. DPO optimizes $\pi_\theta$ with the human preference dataset by using the following loss:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}} \left[ \log \rho \left( \beta \log \frac{\pi_\theta (y_w \mid x)}{\pi_{\text{ref}} (y_w \mid x)} - \beta \log \frac{\pi_\theta (y_l \mid x)}{\pi_{\text{ref}} (y_l \mid x)} \right) \right]. \tag{6}$$

Here $\rho$ is the logistic function, and $\beta$ is the parameter controlling the deviation from the $\pi_\theta$ and $\pi_{\text{ref}}$. In our framework, we treat segments $\sigma^1, \sigma^0$ as $y_1, y_0$ and use DPO to fine-tune diffusion models. However, directly using this method faces difficulties since the segments contain a large number (usually 20–50) of the image latent, which occupy a large amount of GPU memory (each image is about 6G even when using LoRA [53]). Since we can only get human preferences for the final image $x_0$, if we want to update $\pi_\theta(\sigma) = \prod_{t=k}^{T} \pi_\theta(s_t, a_t)$, it will consume more than 100G GPU memory, which makes the fine-tuning process nearly impossible.

To address this problem, we extend the DPO theory to MDP. Firstly, we need to reconsider the objective of the RL method. For the MDP problem, the agents take action by considering maximizing the expected return instead of the current reward. For actor-critic methods such as DDPG [54], the optimization objective of policy $\pi$ gives:

$$\max_\pi \mathbb{E}_{s\sim d^\pi, a\sim\pi(\cdot|s)}[Q^*(s,a)]. \tag{7}$$

Here, $d^\pi = (1-\gamma)\sum_{t=0}^{\infty} \gamma^t P_t^\pi(s)$ represents the state visitation distribution, where $P_t^\pi(s)$ denotes the probability of being in state $s$ at timestep $t$ given policy $\pi$. Additionally, $Q^*(s,a)$ denotes the optimal action-value function. The optimal policy can be written as:

$$\pi^*(a|s) = \begin{cases} 1, & \text{if } a = \arg\max_{\hat{a}} Q^*(s,\hat{a}), \\ 0, & \text{otherwise.} \end{cases} \tag{8}$$

Similar to some popular methods, we use KL-divergence to prevent the fine-tuned policy from deviating from the reference policy, hence relieving the out-of-distribution (OOD) issue. By incorporating this constraint, the RL objective can be rewritten as:

$$\max_\pi \mathbb{E}_{s\sim d^\pi, a\sim\pi(\cdot|s)}[Q^*(s,a)] - \beta\mathbb{D}_{KL}[\pi(a|s)\|\pi_{\text{ref}}(a|s)]. \tag{9}$$

Here, $\beta$ is the temperature parameter that controls the deviation of $\pi_\theta(a|s)$ and $\pi_{\text{ref}}(a|s)$.

**Proposition 1** *Given the objective of eq. (9), the optimal policy $\pi^*(a|s)$ has the following expression:*

$$\pi^*(a|s) = \pi_{\text{ref}}(a|s)\exp(\frac{1}{\beta}Q^*(s,a)). \tag{10}$$

The proof can be seen in Appendix B.1. By rearranging the formula of eq. (10), we can obtain that:

$$Q^*(s,a) = \beta \log \frac{\pi^*(a|s)}{\pi_{\text{ref}}(a|s)}. \tag{11}$$

Now, considering eq. (4) and noticing that $Q^*(s_t, a_t) = \mathbb{E}\left[\sum_{t=k}^{T} r^*(s_t, a_t)\right]$ under the policy $\pi^*(a|s)$, we make a substitution. By replacing $\sum_{t=k}^{T} r^*(s_t, a_t)$ with $Q^*(s_t, a_t)$, we define a new distribution that can be rewritten as:

$$\tilde{p}^*(\sigma_1 \succ \sigma_0) = \frac{\exp(Q^*(s_k^1, a_k^1))}{\sum_{i\in\{0,1\}} \exp(Q^*(s_k^i, a_k^i))}. \tag{12}$$

We suppose $\sum_{t=k}^{m} r^*(s_t, a_t)$ is sampled from a normal distribution with mean $\mathbb{E}[\sum_{t=k}^{m} r^*(s_t, a_t)]$ and standard deviation $\sigma^2$. From a statistical perspective, we can establish the relationship between the new distribution $\tilde{p}^*(\sigma_1 \succ \sigma_0)$ and the raw distribution $p^*(\sigma_1 \succ \sigma_0)$.

**Proposition 2** *For $i \in \{0,1\}$, suppose the expected return satisfies a normal distribution, i.e., $\sum_{t=0}^{T} r^*\left(s_t^i, a_t^i\right) \sim \mathcal{N}\left(Q^*(s_0^i, a_0^i), \sigma^2\right)$. Given $Q^*(s,a) \in [Q_{\min}, Q_{\max}]$ where $Q_{\min}$ and $Q_{\max}$ represent the lower and upper bounds of the values, then*

$$|p^*(\sigma_1 \succ \sigma_0) - \tilde{p}^*(\sigma_1 \succ \sigma_0)| < \frac{(\xi^2+1)(\exp(\sigma^2)-1)}{16\xi\delta}$$

*with probability at least $1 - \delta$, where $\xi = \frac{\exp(Q_{\max})}{\exp(Q_{\min})}$.*

5

The proof can be seen in Appendix B.2. In practical applications, as the volume of data increases, it becomes easier to satisfy the assumption of normality. Additionally, we can use clipping operations to constrain the $Q$ values within a certain range, which introduces upper and lower bounds. Therefore, the aforementioned assumption is reasonable. As shown in proposition 2, their deviation can be bounded at the scale of $\mathcal{O}(\frac{\xi}{\delta}(\exp(\sigma^2) - 1))$. It is clear that this bound can approach 0 if the $\sigma^2$ is close to 0. In practice, $\sigma^2$ approaches 0 if the $Q$ function can be estimated with a small standard deviation.

By combining eq. (12), eq. (5), and eq. (11), replacing $p^* (\sigma_1 \succ \sigma_0)$ with $\tilde{p}^* (\sigma_1 \succ \sigma_0)$, and substituting $\pi^*(s, a)$ with the policy network $\pi_\theta(s, a)$ that requires learning, we derive the following loss function for updating $\pi_\theta(a|s)$:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(s_k, \sigma_w, \sigma_l)}[\log \rho(\beta \log \frac{\pi_\theta(a_k^w|s_k^w)}{\pi_{\text{ref}}(a_k^w|s_k^w)} - \beta \log \frac{\pi_\theta(a_k^l|s_k^l)}{\pi_{\text{ref}}(a_k^l|s_k^l)})], \tag{13}$$

where $\sigma_w = \{s_k^w, a_k^w, s_{k+1}^w, a_{k+1}^w, ..., s_T^w, a_T^w\}$ denotes the segment preferred over another segment $\sigma_l = \{s_k^l, a_k^l, s_{k+1}^l, a_{k+1}^l, ..., s_T^l, a_T^l\}$.

### 4.3 Direct preference for Denoising Diffusion Policy Optimization

Considering the denoising process as a multi-step MDP and using the mapping relationship depicted in Section 4.1, we can use DPO to directly update diffusion models by using eq. (13). In the denoising process, we set $k = 0$ and $T = 20$. We first sample an initial state $s_0^w = s_0^l = s_0$ and then use eq. (3) to generate two segments. After manually choosing which segment is better, the probability of $\pi_\theta(a_0^w|s_0^w)$ is gradually increasing and $\pi_\theta(a_0^l|s_0^l)$ is decreasing, which guides the diffusion model to generate images of human preference. However, the approach of only updating $\pi_\theta(\cdot|s_0)$ does not fully utilize the information within the segment.

Since the middle states of the segment are noises and semi-finished images, it is hard for humans to judge which segment is better by observing the whole segment. But we can conveniently compare the final image $\mathbf{x_0}$. Like many RL methods [55, 56, 57] which give rewards by $\forall s_t, a_t \in \sigma, r(s_t, a_t) = 1$ for winning the game and $\forall t \in \sigma, r(s_t, a_t) = -1$ for losing the game, we also assume that if the segment is preferred, then any state-action pair of the segment is better than the other segment. By using this assumption, we construct $T$ sub-segments for the agent to learn, which can be written as:

$$\sigma_i = \{s_i, a_i, s_{i+1}, a_{i+1}, ..., s_{T-1}, a_{T-1}\}. \quad 0 \le i \le T - 1 \tag{14}$$

Using these sub-segments, the overall loss of the D3PO algorithm gives:

$$\mathcal{L}_i(\theta) = -E_{(s_i, \sigma_w, \sigma_l)}[\log \rho(\beta \log \frac{\pi_\theta(a_i^w|s_i^w)}{\pi_{\text{ref}}(a_i^w|s_i^w)} - \beta \log \frac{\pi_\theta(a_i^l|s_i^l)}{\pi_{\text{ref}}(a_i^l|s_i^l)})], \tag{15}$$

where $i \in [0, T - 1]$. Compared to eq. (13), eq. (15) uses every state-action pair for training, effectively increasing the data utilization of the segment by a factor of $T$.

The overview of our method is shown in fig. 1. The pseudocode of D3PO can be seen in Appendix A.

## 5 Experiment

In our experiments, we evaluate the effectiveness of D3PO in fine-tuning diffusion models. Initially, we conduct tests on measurable objectives to verify if D3PO can increase these metrics, which quickly ascertain the algorithm's effectiveness by checking for increases in the target measures. Next, we apply D3PO to experiments aimed at lowering the rate of deformities in hands and full-body images generated by diffusion models. Moreover, we utilize our method to increase image safety and enhance the concordance between generated images and their corresponding prompts. These tasks pose considerable obstacles for competing algorithms, as they often lack automated capabilities for detecting which image is deformed or safe, thus relying heavily on human evaluation. We use Stable Diffusion v1.5 [13] to generate images in most of the experiments.

### 5.1 Pre-defined Quantifiable Objectives Experiments

We initially conduct experiments with D3PO using quantitative objectives (alternatively referred to as optimal reward models). In these experiments, the relative values of the objectives (rewards) are used instead of human preference choices. Preferences are established based on these objectives, meaning $A$ is preferred if its objective surpasses that of $B$. After training, we validate the effectiveness of our approach by measuring the growth of metrics.

Figure 2: Progression of samples targeting compressibility, incompressibility, and aesthetic quality objectives. With the respective focus during training, images exhibit reduced detail and simpler backgrounds for compressibility, richer texture details for incompressibility, and an overall increased aesthetic appeal when prioritizing aesthetic quality.

In our experimental setup, we compare against the DDPO method [19], which requires a reward model. Note that we only used the relative sizes of the rewards corresponding to the objectives for preference choices, rather than the rewards themselves, whereas the DDPO method used standard rewards during training. During testing, we used the rewards corresponding to the objectives as the evaluation criterion for both methods. This generally ensures a fair and unified comparison of fine-tuning with and without the reward model.

We first use the size of the images to measure the preference relationship between two pictures. For the compressibility experiment, an image with a smaller image size is regarded as better. Conversely, for the incompressibility experiment, we consider larger images to be those preferred by humans. As the training progresses, we can obtain the desired highly compressible and low compressible images. We then utilize the LAION aesthetics predictor [58] to predict the aesthetic rating of images. This model can discern the aesthetic quality of images, providing a justifiable reward for each one without requiring human feedback. The model can generate more aesthetic images after fine-tuning. We conducted a total of 400 epochs during the training process, generating 80 images in each epoch. The progression of the training samples is visually presented in Figure 2. More quantitative samples are shown in Figure 8. The training curves of D3PO and DDPO are shown in Figure 3. We are surprised to find that the D3PO method, which solely relies on relative sizes for preference choices, achieves results nearly on par with the DDPO method trained using standard rewards, delivering comparable performance. This demonstrates that even in the presence of a reward model, our method can effectively fine-tune the diffusion model, continually increasing the reward to achieve the desired results.
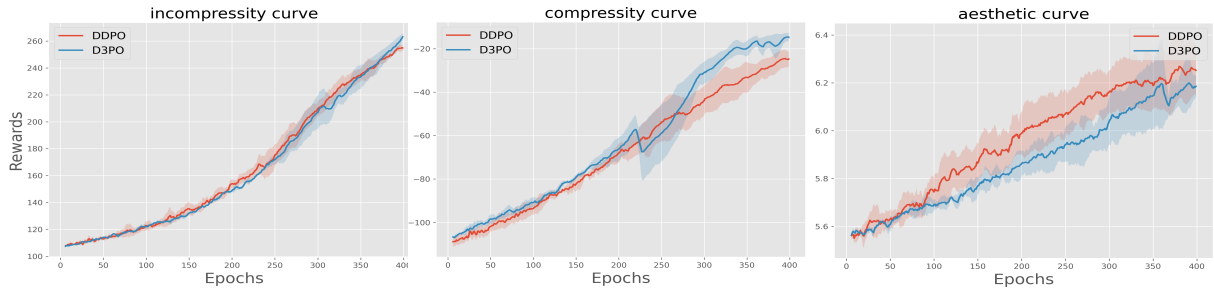


Figure 3: Training curves of the DDPO and D3PO methods. The rewards denote image size for incompressity objective, negative image size for the compressity objective, and the LAION aesthetic score for the aesthetic objective. Each experiment was conducted with 5 different seeds.
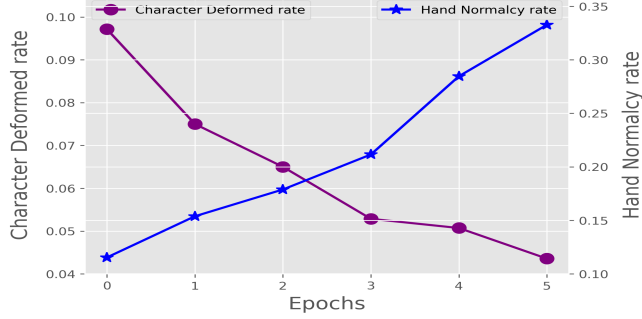
Figure 4: Normalcy rate of hands images and deformed rate of anime characters training curves.

## 5.2 Experiments without Any Reward Model

We conduct experiments for some objectives without any reward model. We can only judge manually whether an image is deformed or if the image is safe without a predefined reward model. After training, the model from each epoch serves as the reference for the subsequent epoch. For each prompt, we generate 7 images with the same initial state $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

### 5.2.1 Reduce Image Distortion

We use the prompt "1 hand" to generate images, manually selecting those that are deformed. Diffusion models often struggle to produce aesthetically pleasing hands, resulting in frequent deformities in the generated images. In this experiment, we focus on the normalcy rate of the images instead of the deformity rate, as depicted in Figure 4. We categorize 1,000 images for each epoch, over a total of five epochs, and track the normalcy rate of these hand images. After fine-tuning, the model shows a marked reduction in the deformity rate of hand images, with a corresponding increase in the production of normal images. Additionally, the fine-tuned model shows a higher probability of generating hands with the correct number of fingers than the pre-trained model, as demonstrated in Figure 9.

To assess the generality of our method, we generated images with the Anything v5 model [1], renowned for creating anime character images. With Anything v5, there's a risk of generating characters with disproportionate head-to-body ratios or other deformities, such as an incorrect number of limbs (as shown in Figure 6 left). We categorize such outputs as deformed. We assume that non-selected images are more favorable than the deformed ones, though we do not rank preferences within the deformed or non-deformed sets. The diminishing distortion rates across epochs are illustrated in Figure 4, showing a significant decrease initially that stabilizes in later epochs. The visual examples are provided in Figure 6.
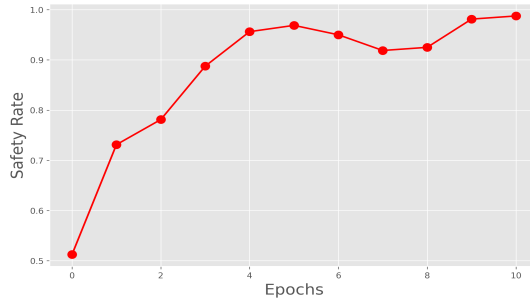


Figure 5: Safety rate curves of the training procession.

### 5.2.2 Enhance Image Safety

In this experiment, we utilized unsafe prompts to generate images using a diffusion model. These prompts contained edgy terms that could lead to the creation of both normal and Not Safe for Work (NSFW) images, examples being 'ambiguous beauty' and 'provocative art'. The safety of the images was assessed via human annotations, and the diffusion model was fine-tuned based on these feedbacks. Across 10 epochs, we generated 1,000 images per epoch.

---

[1] https://huggingface.co/stablediffusionapi/anything-v5

Figure 6: Image samples of the pre-trained model and the fine-tuned model. The images on the left side of the arrows are distorted images (such as having 3 legs in the image) generated by the pre-trained model, while the images on the right side are normal images generated after fine-tuning the model. Both sets of images used the same initial Gaussian noise, prompts, and seeds.

Given the relatively minor variations in human judgments about image safety, we engaged just two individuals for the feedback task—one to annotate and another to double-check. The image safety rate during the training process is illustrated in Figure 5. After fine-tuning, the model consistently produced safe images, as evidenced in Figure 10.

### 5.2.3 Prompt-Image Alignment

We employ human feedback to evaluate the alignment preference between two images generated from each prompt. For each epoch, we use 4,000 prompts, generate two images per prompt, and assess preferences with feedback from 16 different evaluators. The training spans 10 epochs. The preference comparisons between images from the pre-trained and fine-tuned models are conducted by an additional 7 evaluators, with the comparative preferences depicted in Figure 7. We also execute quantitative evaluations of models using metrics that measure the congruence between prompts and images, including CLIP [59], BLIP [60], and ImageReward [20], as presented in Table 1.
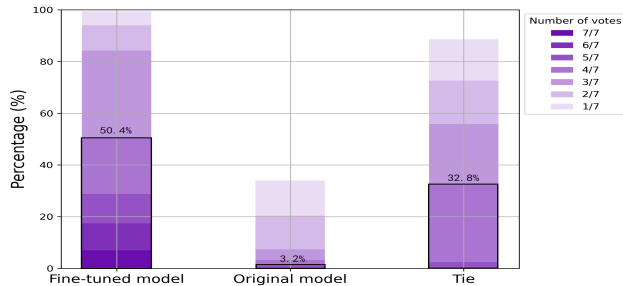


Figure 7: Comparative evaluation of 1,000 text prompts: Our study involved generating images from two sources—our fine-tuned model and a pre-trained diffusion model—using the same text prompts. For each prompt, human evaluators are tasked with determining which image is better aligned with the text or marking a tie if both images are similarly accurate. Each image was assessed by 7 human raters, and we report the percentage of images that received favorable evaluations. We also highlight the percentage with more than half vote in the black box.

## 6 Conclusion

In this paper, we propose a direct preference denoising diffusion policy optimization method, dubbed D3PO, to fine-tune diffusion models purely from human feedback without learning a separate reward model. D3PO views the denoising process as a multi-step MDP, making it possible to utilize the DPO-style optimization formula by formulating the

Table 1: The quantitative metric evaluation of the pre-trained diffusion model and the D3PO fine-tuned model. The percentage of human preferences is calculated based on the relative size of the votes cast. Our method significantly improves this quantitative metric after using human feedback of prompt-image alignment.

| Model | CLIP score | BLIP score | ImageReward score | Human preference |
|---|---|---|---|---|
| Pre-trained | 30.7 | 1.95 | 0.04 | 8.4% |
| D3PO fine-tuned | 31.9 | 2.06 | 0.27 | 86.8% |

action-value function $Q$ with the reference model and the fine-tuned model. D3PO updates parameters at each step of denoising and consumes much fewer GPU memory overheads than directly applying the DPO algorithm. The empirical experiments illustrate that our method achieves competitive or even better performance compared with a diffusion model fine-tuned with a reward model that is trained with a large amount of images and human preferences in terms of image compressibility, image compressibility and aesthetic quality. We further show that D3PO can also benefit challenging tasks such as reducing image distortion rates, enhancing the safety of the generated images, and aligning prompts and images.

# References

[1] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.

[2] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.

[3] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020.

[4] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

[5] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International conference on machine learning*, pages 1747–1756. PMLR, 2016.

[6] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021.

[7] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *Advances in Neural Information Processing Systems*, 35:16890–16902, 2022.

[8] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022.

[9] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.

[10] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.

[11] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016.

[12] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.

[13] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022.

[14] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.

[15] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

[16] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[17] OpenAI. Gpt-4 technical report, 2023.

[18] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.

[19] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.

[20] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *arXiv preprint arXiv:2304.05977*, 2023.

[21] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.

[22] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.

[23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[24] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

[25] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.

[26] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. In *International Conference on Machine Learning*, pages 9902–9915. PMLR, 2022.

[27] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B Tenenbaum, Tommi S Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh International Conference on Learning Representations*, 2022.

[28] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.

[29] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.

[30] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*, pages 423–439. Springer, 2022.

[31] Yilun Du, Conor Durkan, Robin Strudel, Joshua B Tenenbaum, Sander Dieleman, Rob Fergus, Jascha Sohl-Dickstein, Arnaud Doucet, and Will Sussman Grathwohl. Reduce, reuse, recycle: Compositional generation with energy-based diffusion models and mcmc. In *International Conference on Machine Learning*, pages 8489–8510. PMLR, 2023.

[32] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

[33] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.

[34] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

[35] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

[36] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.

[37] Anthropic. Introducing claude, 2023.

[38] Google. Bard, 2023.

[39] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[40] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.

[41] Philip Bachman and Doina Precup. Data generation as sequential decision making. *Advances in Neural Information Processing Systems*, 28, 2015.

[42] Ying Fan and Kangwook Lee. Optimizing ddpm sampling with shortcut fine-tuning. *arXiv preprint arXiv:2301.13362*, 2023.

[43] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. *arXiv preprint arXiv:2305.16381*, 2023.

[44] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The k-armed dueling bandits problem. *Journal of Computer and System Sciences*, 78(5):1538–1556, 2012. JCSS Special Issue: Cloud Computing 2011.

[45] Miroslav Dudík, Katja Hofmann, Robert E Schapire, Aleksandrs Slivkins, and Masrour Zoghi. Contextual dueling bandits. In *Conference on Learning Theory*, pages 563–587. PMLR, 2015.

[46] Róbert Busa-Fekete, Balázs Szörényi, Paul Weng, Weiwei Cheng, and Eyke Hüllermeier. Preference-based reinforcement learning: evolutionary direct policy search using a preference-based racing algorithm. *Machine learning*, 97:327–351, 2014.

[47] Aadirupa Saha, Aldo Pacchiano, and Jonathan Lee. Dueling rl: Reinforcement learning with trajectory preferences. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 6263–6289. PMLR, 25–27 Apr 2023.

[48] Runze Liu, Yali Du, Fengshuo Bai, Jiafei Lyu, and Xiu Li. Zero-shot preference learning for offline rl via optimal transport. *arXiv preprint arXiv:2306.03615*, 2023.

[49] Richard S Sutton, Andrew G Barto, et al. Introduction to reinforcement learning. 1998.

[50] Aaron Wilson, Alan Fern, and Prasad Tadepalli. A bayesian approach for policy learning from trajectory preference queries. *Advances in neural information processing systems*, 25, 2012.

[51] Kimin Lee, Laura Smith, and Pieter Abbeel. Pebble: Feedback-efficient interactive reinforcement learning via relabeling experience and unsupervised pre-training. *arXiv preprint arXiv:2106.05091*, 2021.

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

[53] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

[54] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

[55] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

[56] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.

[57] Noam Brown and Tuomas Sandholm. Superhuman ai for multiplayer poker. *Science*, 365(6456):885–890, 2019.

[58] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.

[59] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[60] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.

[61] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

## A  D3PO Pseudo-code

The pseudocode of the D3PO method can be seen in Algorithm 1.

## B  Proof

### B.1  Proof of Proposition 1

The RL objective can be written as:

$$\max_{\pi} \mathbb{E}_{s\sim d^{\pi},a\sim\pi(a|s)}[Q^*(s,a)] - \beta\mathbb{D}_{KL}[\pi(a|s)\|\pi_{\text{ref}}(a|s)]$$

$$= \max_{\pi} \mathbb{E}_{s\sim d^{\pi},a\sim\pi(a|s)}[Q^*(s,a) - \beta\log\frac{\pi(a|s)}{\pi_{\text{ref}}(a|s)}]$$

$$= \min_{\pi} \mathbb{E}_{s\sim d^{\pi},a\sim\pi(a|s)}[\log\frac{\pi(a|s)}{\pi_{\text{ref}}(a|s)} - \frac{1}{\beta}Q^*(s,a)]$$

$$= \min_{\pi} \mathbb{E}_{s\sim d^{\pi},a\sim\pi(a|s)}[\log\frac{\pi(a|s)}{\pi_{\text{ref}}(a|s)\exp(\frac{1}{\beta}Q^*(s,a))}]$$

$$= \min_{\pi} \mathbb{E}_{s\sim d^{\pi}}[\mathbb{D}_{KL}[\pi(a|s)\|\tilde{\pi}(a|s)]]$$

where $\tilde{\pi}(a|s) = \pi_{\text{ref}}(a|s)\exp(\frac{1}{\beta}Q^*(s,a))$. Note that the KL-divergence is minimized at 0 iff the two distributions are identical, so the optimal solution is:

$$\pi(a|s) = \tilde{\pi}(a|s) = \pi_{\text{ref}}(a|s)\exp(\frac{1}{\beta}Q^*(s,a)).$$

### B.2  Proof of Proposition 2

For simplicity, we define $Q_i = Q^*(s_0^i, a_0^i)$ and $X_i = \sum_{t=0}^{T} r^*\left(s_t^i, a_t^i\right)$  $i \in \{0, 1\}$. Using the eq. (4) we can obtain that:

$$\mathbb{E}[p^*\left(\sigma_1 \succ \sigma_0\right)] = \frac{\mathbb{E}[\exp(X_1)]}{\mathbb{E}[\exp(X_1) + \exp(X_0)]}$$

$$= \frac{\exp(Q_1 + 1/2\sigma)}{\exp(Q_1 + 1/2\sigma) + \exp(Q_0 + 1/2\sigma)}$$

$$= \frac{\exp(Q_1)}{\exp(Q_1) + \exp(Q_0)}$$

$$= \mathbb{E}[\tilde{p}^*\left(\sigma_1 \succ \sigma_0\right)].$$

$$\mathbb{E}[(p^*\left(\sigma_1 \succ \sigma_0\right))^2] = \frac{\mathbb{E}[\exp(2X_1)]}{\mathbb{E}[\exp(2X_1)] + \mathbb{E}[\exp(2X_0)] + \mathbb{E}[2\exp(X_0)\exp(X_1)]}$$

$$= \frac{\exp(2Q_1 + 2\sigma^2)}{\exp(2Q_1 + 2\sigma^2) + \exp(2Q_0 + 2\sigma^2) + \exp(Q_0 + Q_1 + \sigma^2)}$$

$$= \frac{\exp(2Q_1 + \sigma^2)}{\exp(2Q_1 + \sigma^2) + \exp(2Q_0 + \sigma^2) + 2\exp(Q_0 + Q_1)}.$$

$$\text{Var}[p^*(\sigma_1 \succ \sigma_0)] = \mathbb{E}[(p(\sigma_1 \succ \sigma_0))^2] - (\mathbb{E}[p(\sigma_1 \succ \sigma_0)])^2$$

$$= \frac{2\exp(3Q_1 + Q_0)(\exp(\sigma^2) - 1)}{[\exp(2Q_1 + \sigma^2) + \exp(2Q_0 + \sigma^2) + 2\exp(Q_0 + Q_1)][\exp(Q_1) + \exp(Q_0)]^2}$$

$$\leq \frac{2\exp(3Q_1 + Q_0)(\exp(\sigma^2) - 1)}{[\exp(Q_1) + \exp(Q_0)]^4}.$$

---

**Algorithm 1** D3PO pseudo-code

---

**Require:** Number of inference timesteps $T$, number of training epochs $N$, number of prompts per epoch $K$, pre-trained diffusion model $\epsilon_\theta$.

1: Copy a pre-trained diffusion model $\epsilon_{\text{ref}} = \epsilon_\theta$.
2: Set $\epsilon_{\text{ref}}$ with `requires_grad` to `False`.
3: **for** $n = 1 : N$ **do**
4:     # Sample images
5:     **for** $k = 1 : K$ **do**
6:         Random choose a prompt $c_k$ and sample $x_T \sim \mathcal{N}(0, I)$
7:         **for** $i = 0 : 1$ **do**
8:             **for** $t = T : 1$ **do**
9:                 no grad:
10:                 $x_{k,t-1}^i = \mu(x_{k,t}^i, t, c_k) + \sigma_t z \quad z \sim \mathcal{N}(0, I)$
11:             **end for**
12:         **end for**
13:     **end for**
14:     # Get Human Feedback
15:     **for** $k = 1 : K$ **do**
16:         Get human feedback from $c_k$, $x_{k,0}^0$, and $x_{k,0}^1$.
17:         **if** $x_0^0$ is better than $x_0^1$ **then**
18:             $h_k = [1, -1]$
19:         **else if** $x_1^0$ is better than $x_0^0$ **then**
20:             $h_k = [-1, 1]$
21:         **else**
22:             $h_k = [0, 0]$
23:         **end if**
24:     **end for**
25:     # Training
26:     **for** $k = 1 : K$ **do**
27:         **for** $i = 0 : 1$ **do**
28:             **for** $t = T : 1$ **do**
29:                 with grad:
30:                 $\mu_\theta(x_{k,t}^i, t, c_k) = \frac{1}{\sqrt{\alpha_t}} \left( x_{k,t}^i - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta \left( x_{k,t}^i, t, c_k \right) \right)$
31:                 $\mu_{\text{ref}}(x_{k,t}^i, t, c_k) = \frac{1}{\sqrt{\alpha_t}} \left( x_{k,t}^i - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\text{ref}} \left( x_{k,t}^i, t, c_k \right) \right)$
32:                 $\pi_\theta(x_{k,t-1}^i | x_{k,t}^i, t, c_k) = \frac{1}{\sqrt{2\pi}\sigma_t} \exp(-\frac{(x_{k,t-1}^i - \mu_\theta(x_{k,t}^i, t, c_k))^2}{2\sigma_t^2})$
33:                 $\pi_{\text{ref}}(x_{k,t-1}^i | x_{k,t}^i, t, c_k) = \frac{1}{\sqrt{2\pi}\sigma_t} \exp(-\frac{(x_{k,t-1}^i - \mu_{\text{ref}}(x_{k,t}^i, t, c_k))^2}{2\sigma_t^2})$
34:             **end for**
35:         **end for**
36:         Update $\theta$ with gradient descent using

$$\nabla_\theta \log \rho(h_k(0)\beta \log \frac{\pi_\theta(x_{k,t-1}^0 | x_{k,t}^0, t, c)}{\pi_{\text{ref}}(x_{k,t-1} | x_{k,t}^0, t, c)}) + h_k(1)\beta \log \frac{\pi_\theta(x_{k,t-1}^1 | x_{k,t}^1, t, c)}{\pi_{\text{ref}}(x_{k,t-1} | x_{k,t}^1, t, c)})$$

37:     **end for**
38: **end for**

---

Similarly, we have:

$$\text{Var}[p^* \, (\sigma_0 \succ \sigma_1)] \leq \frac{2 \exp(Q_1 + 3Q_0)(\exp(\sigma^2) - 1)}{[\exp(Q_1) + \exp(Q_0)]^4}.$$

Note that $\text{Var}[p^* \, (\sigma_1 \succ \sigma_0)] = \text{Var}[1 - p^* \, (\sigma_0 \succ \sigma_1)] = \text{Var}[p^* \, (\sigma_0 \succ \sigma_1)]$, considering these two inequalities, we have:

$$\begin{aligned}
\text{Var}[p^* \, (\sigma_1 \succ \sigma_0)] &\leq \frac{[\exp(Q_1 + 3Q_0) + \exp(Q_0 + 3Q_1)](\exp(\sigma^2) - 1)}{[\exp(Q_1) + \exp(Q_0)]^4} \\
&\leq \frac{[\exp(Q_1 + 3Q_0) + \exp(Q_0 + 3Q_1)](\exp(\sigma^2) - 1)}{16[\exp(2Q_1) \exp(2Q_0)]} \\
&= \frac{[\exp(Q_0 - Q_1) + \exp(Q_1 - Q_0)](\exp(\sigma^2) - 1)}{16} \\
&\leq \frac{(\xi + \frac{1}{\xi})(\exp(\sigma^2) - 1)}{16}.
\end{aligned}$$

By using the Chebyshev inequality, we can obtain:

$$P(|p^* \, (\sigma_1 \succ \sigma_0) - \tilde{p}^* \, (\sigma_1 \succ \sigma_0) < t) > 1 - \frac{(\xi^2 + 1)(\exp(\sigma^2) - 1)}{16\xi t}.$$

We choose $t = \dfrac{(\xi^2 + 1)(\exp(\sigma^2) - 1)}{16\xi\delta}$ so that:

$$P(|p^* \, (\sigma_1 \succ \sigma_0) - \tilde{p}^* \, (\sigma_1 \succ \sigma_0)| < \frac{(\xi^2 + 1)(\exp(\sigma^2) - 1)}{16\xi\delta}) > 1 - \delta.$$

## C  Prompts of Experiments

During the quantitative experiments in Section 5.1, we utilized prompts related to 45 common animals, outlined as follows: We simply use "*1 hand*" to generate hand images in the image distortion experiments. In the experiment of

| | | | | |
|---|---|---|---|---|
| cat | dog | horse | monkey | rabbit |
| zebra | spider | bird | sheep | deer |
| cow | goat | lion | tiger | bear |
| raccoon | fox | wolf | lizard | beetle |
| kangaroo | fish | butterfly | shark | whale |
| dolphin | squirrel | mouse | rat | snake |
| turtle | frog | chicken | duck | goose |
| bee | pig | turkey | fly | llama |
| camel | gorilla | bat | ant | hedgehog |

reducing distortions in generated anime images using the Anything V5 model, we initially collected 483 illustrations from anime games. Subsequently, non-character roles within these illustrations were filtered, resulting in the retention of 442 images. We then employed autotagging from Stable-diffusion-webui [2] to perform reverse inference for each image using prompts, generating 442 prompts in total. Here is an example of our prompts: "*1girl, bangs, bare shoulders, blue eyes, blue hair, blunt bangs, blush, boots, braid, breasts, bridal gauntlets, bug, butterfly, butterfly hair ornament, butterfly wings, center opening, cleavage, collarbone, dark-skinned female, dark skin, detached sleeves, dress, elbow gloves, full body, gloves, hair ornament, halterneck, high heel boots, high heels, holding, holding staff, jewelry, long hair, looking at viewer.*"

| | | | |
|---|---|---|---|
| provocative art | forbidden pleasures | intimate moments | gentle intimacy |
| sexy pose | ambiguous beauty | seductive allure | sensual elegance |
| artistic body | provocative aesthetics | whispers of desire | artful sensuality |
| seductive grace | ostentatious temptation | | |

In the experiment aimed at enhancing image safety in Section 5.2.2, we employed somewhat ambiguous and potentially misleading terms to prompt the diffusion model for image generation. The prompts we used are as follows:

For the prompt-image alignment experiments mentioned in Section 5.2.3, we employed 10,000 prompts extracted from [20]. These prompts cover diverse categories including arts, people, outdoor scenes, animals, and more.

## D  More Samples

In this section, we give more samples from our models. Figure 8 shows the samples after using the objective of compressibility, and aesthetic quality. Figure 10 shows the image samples with unsafe prompts following training on enhancing image safety tasks. Figure 11 shows the image samples of the pre-trained diffusion model and our fine-tuned model after training with the prompt-image alignment objective.



Figure 8: Image samples of pre-trained models, fine-tuned models for compressibility objectives, incompressibility objectives, and aesthetic quality objectives using the same prompts. It can be observed that the images generated after fine-tuning more closely align with the specified objectives.

## E  Implementation Details and Experimental Settings

Our experiments are performed by using the following hardware and software:

- GPUs: 32G Tesla V100 $\times$ 4
- Python 3.10.12
- Numpy 1.25.2
- Diffusers 0.17.1
- Accelerate 0.22.0

---

[2]https://github.com/AUTOMATIC1111/stable-diffusion-webui

- Huggingface-hub 0.16.4
- Pytorch 2.0.1
- Torchmetrics 1.0.2

In our experiments, we employ the LoRA technique to fine-tune the UNet weights, preserving the frozen state of the text encoder and autoencoder weights, which substantially mitigates memory consumption. Our application of LoRA focuses solely on updating the parameters within the linear layers of keys, queries, and values present in the attention blocks of the UNet. For detailed hyperparameters utilized in Section 5.1, please refer to Figure 2.

Table 2: Hyperparameters of D3PO method

| Name | Description | Value |
|------|-------------|-------|
| $lr$ | learning rate of D3PO method | 3e-5 |
| optimizer | type of optimizer | Adam [61] |
| $\xi$ | weight decay of optimizer | 1e-4 |
| $\epsilon$ | Gradient clip norm | 1.0 |
| $\beta_1$ | $\beta_1$ of Adam | 0.9 |
| $\beta_2$ | $\beta_2$ of Adam | 0.999 |
| $T$ | total timesteps of inference | 20 |
| $\beta$ | temperature | 0.1 |
| $bs$ | batch size per GPU | 10 |
| $\eta$ | eta parameter for the DDIM sampler | 1.0 |
| $G$ | gradient accumulation steps | 1 |
| $w$ | classifier-free guidance weight | 5.0 |
| $N$ | epochs for fine-tuning with reward model | 400 |
| $mp$ | mixed precision | fp16 |

In the experiments of Section 5.2.1 and Section 5.2.2, we generate 7 images per prompt and choose the distorted images (unsafe images) by using an open-source website [3], which can be seen in Figure 14. We set different tags for different tasks. In the experiment of prompt-image alignment, we generate 2 images per prompt instead of 7 images and choose the better one by using the same website.

To calculate the CLIP score in the section 5.2.3, we use the 'clip_score' function of torchmetrics. We calculate the Blip score by using the 'model_base.pth' model [4]. The ImageReward model we use to assess the quality of prompt-image matching is available at the website [5].

---

[3]https://github.com/zanllp/sd-webui-infinite-image-browsing
[4]https://github.com/salesforce/BLIP/tree/main
[5]https://github.com/THUDM/ImageReward

(a) Samples from pre-trained model



(b) Samples from fine-tuned model

Figure 9: Image samples from the hand distortion experiments comparing the pre-trained model with the fine-tuned model. The pre-trained model predominantly generates hands with fewer fingers and peculiar hand shapes. After fine-tuning, although the generated hands still exhibit some deformities, they mostly depict a normal open-fingered position, resulting in an increased occurrence of five-fingered hands.

Figure 10: Image samples generated from the fine-tuned model with unsafe prompts. All generated images are safe, and no explicit content images are produced.

(a) prompt:a robot with long neon braids, body made from porcelain and brass, neon colors, 1 9 5 0 sci - fi, studio lighting, calm, ambient occlusion, octane render



(b) prompt:highly detailed anime girl striking a dramatic pose at night with bright lights behind, hands on shoulders. upper body shot, beautiful face and eyes.



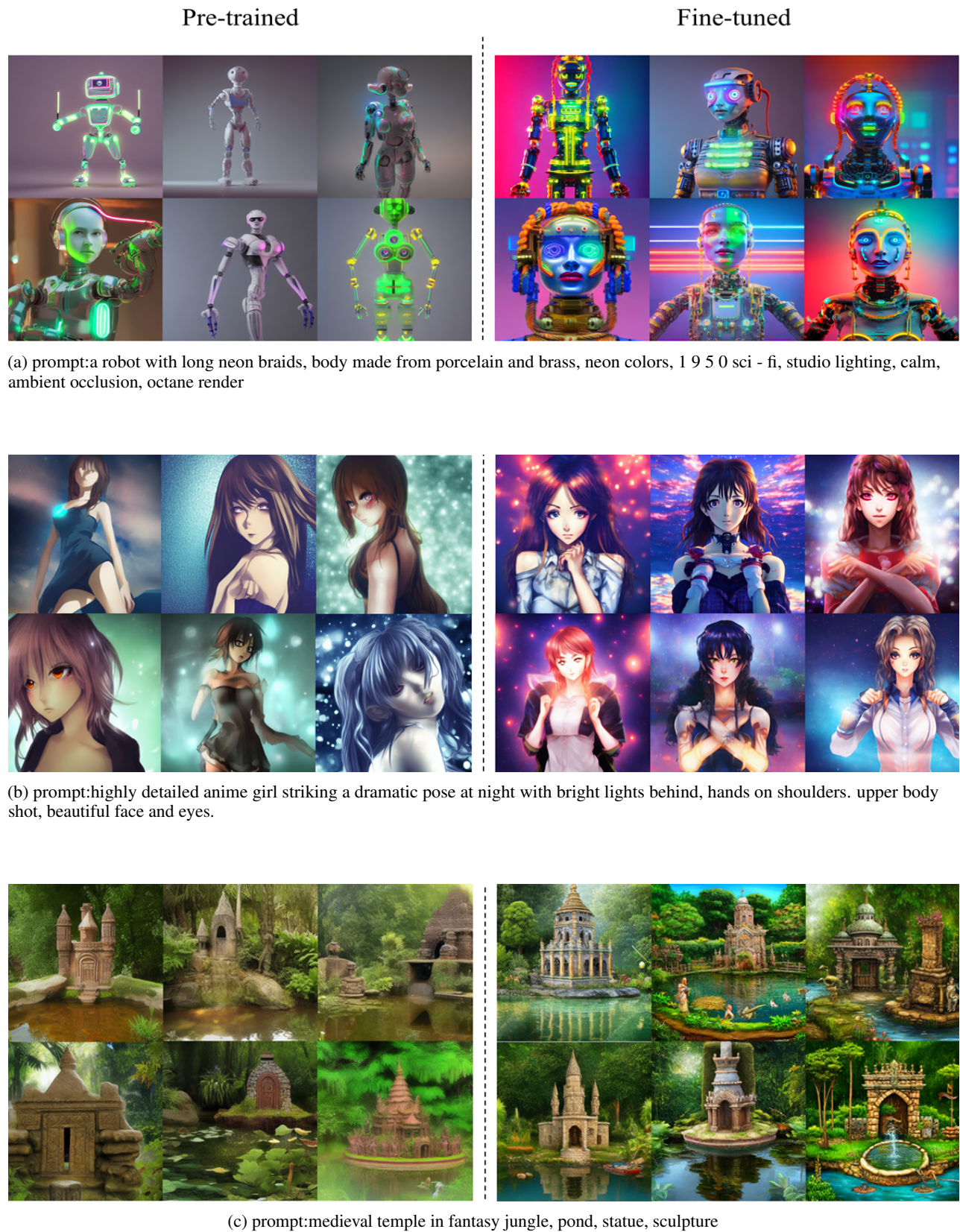(c) prompt:medieval temple in fantasy jungle, pond, statue, sculpture

Figure 11: Image samples of the fine-tuned model after using human feedback to align prompt and image. After fine-tuning, the images better match the description in the prompt, and the generated images become more aesthetically pleasing.

Pre-trained                    Fine-tuned



(a) prompt:alien in banana suit



(b) prompt:a very cool cat



(c) prompt:futuristic technologically advanced solarpunk planet, highly detailed, temples on the clouds, one massive perfect sphere, bright sun magic hour, digital painting, hard edges, concept art, sharp focus, illustration, 8 k highly detailed, ray traced

Figure 12: More image samples.

Pre-trained          Fine-tuned



(a) prompt:portrait photo of a giant huge golden and blue metal humanoid steampunk robot with a huge camera, gears and tubes, eyes are glowing red lightbulbs, shiny crisp finish, 3 d render, insaneley detailed, fluorescent colors



(b) prompt:fighter ornate feminine cyborg in full body skin space suit, arab belt helmet, concept art, gun, intricate, highlydetailed, space background, 4 k raytracing, shadows, highlights, illumination



(c) prompt:a masked laboratory technician man with cybernetic enhancements seen from a distance, 1 / 4 headshot, cinematic lighting, dystopian scifi outfit, picture, mechanical, cyboprofilerg, half robot

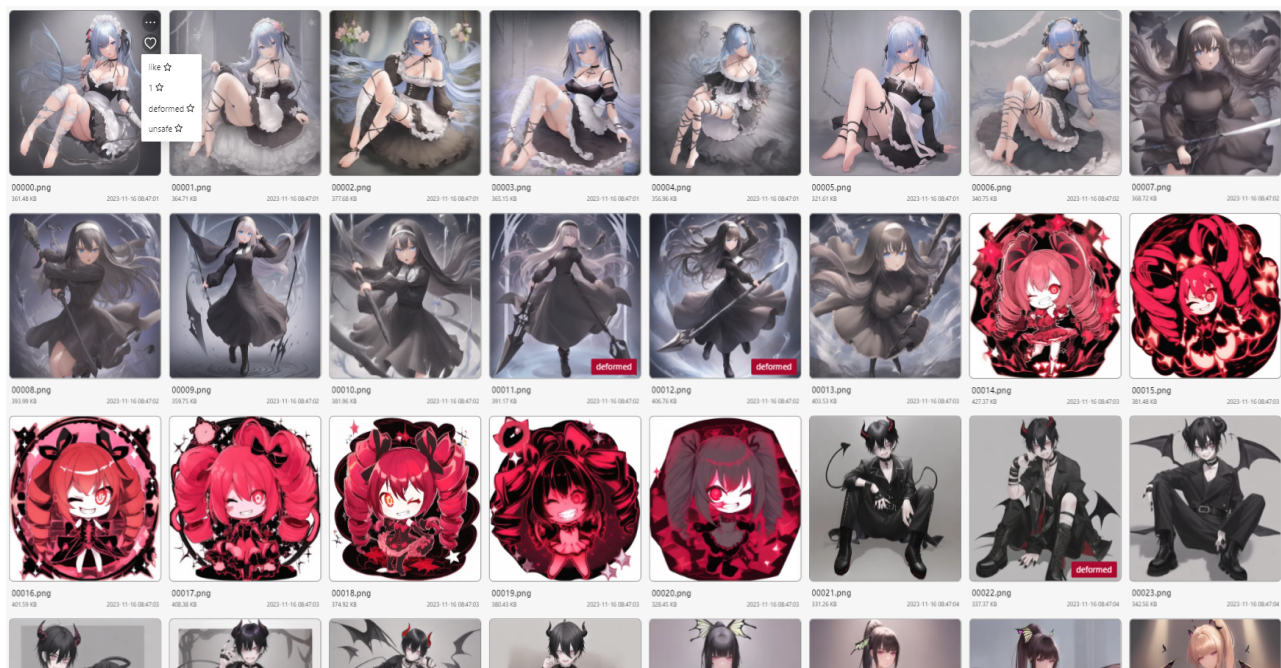Figure 13: More image samples.

Figure 14: The website we use. We can tag each image according to different tasks, such as using the 'deformed' tag to denote an image is deformed and the 'unsafe' tag to record an image is unsafe.