
Improvement-Guided Iterative DPO for Diffusion Models

Ying Fan¹ Fei Deng² Yang Zhao³ Sahil Singla³ Rahul Jain³ Tingbo Hou³ Kangwook Lee¹ Feng Yang³
Deepak Ramachandran³ Qifei Wang³

Abstract

Direct Preference Optimization (DPO) has been shown to be an effective solution in aligning generative models with human preferences. The recent deep dive shows that DPO’s performance is constrained by the offline preference dataset. To solve this challenge, this paper introduces a novel improvement-guided approach for online iterative optimization of the diffusion models without extra annotation. We propose to learn an improvement model to extract the implicit preference improvement direction from the preference dataset. The learned improvement model is then used to generate winning images given the images generated by the current diffusion model as losing images. Thus, the improvement model can guide iterative DPO by generating such online preference datasets repeatedly. This method enables online improvement beyond offline DPO training without requiring additional human labeling or risking overfitting the reward model. Results demonstrate improvements in preference alignment with higher diversity compared with other fine-tuning methods. Our work bridges the gap between offline preference learning and online improvement, offering a promising direction for enhancing diffusion models in image generation tasks with limited preference data.

1. Introduction

Reinforcement Learning from Human Feedback (RLHF) has emerged as a powerful paradigm for aligning generative models with human preferences, showing remarkable success in both language models (Ouyang et al., 2022) and diffusion models for image generation (Black et al., 2024; Fan et al., 2023). Traditional RLHF approaches, often implemented using Proximal Policy Optimization (PPO) (Schulman et al., 2017), have faced significant challenges,

including training instability, overfitting the reward model, and high computational costs.

Direct Preference Optimization (DPO) was introduced as an alternative that simplifies the training process by directly optimizing the model based on preference data (Rafailov et al., 2024; Wallace et al., 2024). While DPO offers more efficient training and improved stability, it is inherently limited to offline dataset, potentially constraining its performance: Offline DPO, which relies solely on a fixed dataset of preferences, often exhibits suboptimal performance due to the lack of on-policy data (Xu et al., 2024; Tajwar et al., 2024). Recent studies have investigated online iterative DPO methods, such as online annotated preference data from LLMs (Rosset et al., 2024), reward models (Xu et al., 2024), or human feedbacks (Xiong et al., 2024). However, online labeling can be prohibitively expensive, and runs the risk of reward hacking (Zhang et al., 2024).

On the other hand, recent research has explored the concept of **self-improvement** in generative models especially LLMs, including self-rewarding models (Yuan et al., 2024b) and self-improving language models (Choi et al., 2024). The core idea of these approaches is to provide reward signals or improvement guidance from some pre-trained model to guide the iterative training process, which offers a promising direction to achieve self-improvement without extra labeled data. As a result, these approaches could be a natural remedy for data constraints in DPO.

The self-improvement approaches have been widely studied for LLMs since the base LLM can be re-purposed in a natural way to provide a self-improvement signal. However, such self-improvement capability remains under-explored in text-to-image diffusion models since it is not straightforward to directly apply the self-improvement approach to a mixed-modality T2I model, which is the main question we focus on in this paper. Recently, Yuan et al. (2024a) proposed a self-play approach for diffusion models. However, their optimization target is equivalent to aligning with the winning data distribution, so the performance is thus still upper-bounded by the offline dataset¹.

In this paper, we aim to answer the following research question: *Can we achieve iterative improvement for diffusion*

^{*}Equal contribution ¹Department of XXX, University of Wisconsin-Madison ²Google DeepMind ³Google. Correspondence to: Qifei Wang <qfwang@google.com>.

models without extra annotations?

To solve the challenges of limited offline annotations, we introduce a novel method to train a *improvement* model to learn the generic *improvement directions* from large-scale preference datasets. Then we use the improvement diffusion model to generate a winning image given a on-policy losing image and a prompt. By iterative DPO with such guidance, the improvement model could exploit the improving guidance from the offline dataset. By extrapolating such knowledge in iterative DPO, the online sampled images can be further improved and potentially break the constraints introduced by the offline dataset. This approach offers several key advantages: Leveraging the advantages of DPO while mitigating its limitations in offline settings; Enabling online learning without extra annotations; Providing a mechanism for iterative improvement for diffusion models with fixed preference datasets.

Experimental results demonstrate that our iterative training for diffusion models guided by learned improvement model leads to improvements over DPO baselines including Diffusion-DPO (Wallace et al., 2024) and SPIN (Yuan et al., 2024a). Specifically, we observe consistently higher scores on PickScore (Kirstain et al., 2023), HPSv2 (Wu et al., 2023), and Aesthetic score (Schuhmann et al., 2022), indicating improved image quality and better alignment with human preferences.

We summarize our contributions as follows: **(1)** We introduce a novel improvement diffusion model that learns an improvement direction from an offline preference dataset. **(2)** Using the improvement model to generate online training data, we address the critical challenge of learning from limited offline preference data, enabling iterative improvement during training. **(3)** Our experiments demonstrate improvements in preference alignment and visual quality compared with baseline DPO methods.

2. Preliminaries

2.1. Diffusion Models

Let $x_0 \in \mathbb{R}^n$ be a data sample, and q_0 be the data distribution, i.e., $x_0 \sim q_0(x_0)$. Diffusion models approximate q_0 with $p_\theta(x_0) = \int p_\theta(x_{0:T}) dx_{1:T}$, where $p_\theta(x_{0:T}) = p_T(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$ is a Markov chain with the following dynamics:

$$p(x_T) = \mathcal{N}(0, I), \quad (1)$$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t, t), \Sigma_t). \quad (2)$$

The *forward* or *diffusion process* $q(x_{1:T}|x_0)$ is a Markov chain that adds Gaussian noise to the data according to a

variance schedule β_1, \dots, β_T :

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad (3)$$

$$q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{1-\beta_t} x_{t-1}, \beta_t I). \quad (4)$$

Let $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$, $\tilde{\beta}_t = \frac{1-\bar{\alpha}_t-1}{1-\bar{\alpha}_t} \beta_t$, $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$. The training of diffusion models is performed by optimizing a variational bound on the negative log-likelihood $\mathbb{E}_q[-\log p_\theta(x_0)]$, which is equivalent to optimizing:

$$\mathbb{E}_{x_t, \epsilon, t} \left[\|\epsilon - \epsilon_\theta(x_t, t)\|_2^2 \right], \quad (5)$$

where $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon$, $x_0 \sim q_0(x_0)$, $\epsilon \sim \mathcal{N}(0, I)$.

2.2. DPO and Diffusion-DPO

DPO. Assume that we have access to a general preference dataset $\mathcal{D} = \{c, x_w, x_l\}$ where c is the text prompt, x_l is the losing response and x_w is the winning response. Given a conditional generative model $p_\theta(x|c)$ and a reference model $p_{\text{ref}}(x|c)$, we can align the model with the preference using the DPO loss (Rafailov et al., 2024):

$$-\mathbb{E}_{c, x_w, x_l \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{p_\theta(x_w|c)}{p_{\text{ref}}(x_w|c)} - \beta \log \frac{p_\theta(x_l|c)}{p_{\text{ref}}(x_l|c)} \right) \right]. \quad (6)$$

Diffusion-DPO. For diffusion models, since $p_\theta(x|c)$ is not generally tractable, (Wallace et al., 2024) proposes an approximation by finding an upper-bound of the original DPO objective:

$$\begin{aligned} & -\mathbb{E}_{c, x_l, x_w \sim \mathcal{D}, t} [\log \sigma (-\beta T (\|\epsilon^w - \epsilon_\theta(x_t^w, t, c)\|_2^2 \\ & \quad + \|\epsilon^l - \epsilon_{\text{ref}}(x_t^l, t, c)\|_2^2 \\ & \quad - \|\epsilon^w - \epsilon_{\text{ref}}(x_t^w, t, c)\|_2^2 \\ & \quad - \|\epsilon^l - \epsilon_\theta(x_t^l, t, c)\|_2^2))] \end{aligned} \quad (7)$$

Drawbacks of DPO. The interpretation of DPO training is straightforward: it aims to pull up the probability of the winning response and pull down the losing one. During training, all the responses are from the preference dataset, and the actual output of the model is never checked. The quality of the learned policy in DPO can be compromised by a biased distribution towards unseen responses. This bias arises when the offline preference dataset lacks diversity or is not readily accessible. This phenomenon has been observed in (Xu et al., 2024).

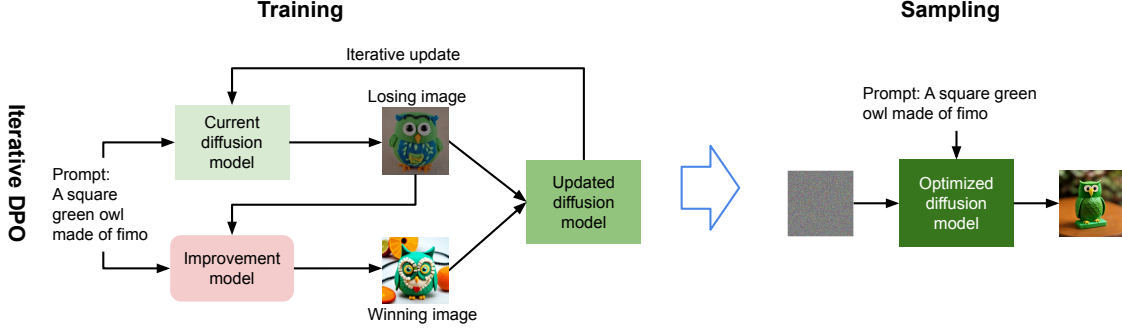


Figure 1. The overview of the pipeline of the improvement-guided iterative DPO process. The diagram on the left side demonstrates the iterative DPO algorithm guided by the improvement model. The current diffusion model generates a losing image and passes it to the improvement model to improve to a winning image. Both images are paired as the preference dataset to fine-tune the diffusion model. The diffusion model is optimized iteratively until it converges. The optimized diffusion model is then deployed for inference as shown in the diagram on the right side.

Given the downsides of using an offline dataset in DPO, the recent work (Xiong et al., 2024) has explored augmenting training datasets through online training, incorporating online samples that enhance performance in preference learning (Tajwar et al., 2024). However, annotating these samples requires extra effort, and optimizing with a reward model could risk reward over-optimization and hacking. This paper explores whether DPO-based training without extra annotations can be further improved.

3. Method

We consider a scenario where *only a fixed offline preference dataset is available, without access to additional annotation sources*. We propose to build an *improvement model* from the preference dataset that generates improved images (for a given prompt) when given images generated by the current diffusion model as input. The input (image condition) and output (improved image) of the improvement model therefore correspond to a losing/winning preference pair that can be used for iterative DPO training without extra annotation.

The intuition behind iterative DPO training with an improvement model is straightforward. Recall that DPO training pulls up the probability of the winning response, which is the output of the improvement model in our case. It also pulls down the probability of the losing response, which is the output of the current diffusion model. Thus, if we can successfully train an improvement model, we can continuously improve the current diffusion model using the improvement model with iterative DPO training till convergence.

We introduce how to train such an improvement model in Section 3.1, the sampling from the improvement model in Section 3.2, and iterative training of the improvement model in Section 3.3.

3.1. Training the Improvement Diffusion Model

The objective of the improvement diffusion model ϕ is to predict a conditional distribution over improved images, $p_{\phi}^{\dagger}(x_w|x_l, c)$ i.e. It learns to generate a winning image x_w given a text prompt c and a losing image x_l . This can be accomplished by the ability of diffusion models to condition the denoising trajectory on arbitrary additional signals. For example, InstructPix2Pix (Brooks et al., 2023) is an image-editing model that takes an original image and an editing instruction (in text) as its input conditions by encoding the image with additional channels in the first convolutional layer of the UNet.

Multi-task training. To train a model with a generic improvement capability to map any given image to higher quality ones without sacrificing diversity, we design a multi-task training algorithm that takes different text-image condition combinations as input (See the left side of Figure 2). The model is trained on a mixture of the following tasks:

1. Learning the conditional winning distribution: Given both text c and a losing image condition x_l , we learn the target distribution of x_w :

$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), x_w, x_l, c, t} [\|\epsilon - \epsilon_{\phi}(x_t|x_l, c, t)\|^2], \quad (8)$$

where $x_t = \sqrt{\alpha_t}x_w + \sqrt{1 - \alpha_t}\epsilon$.

2. Reconstruction: Given only the image condition $x \in \{x_w, x_l\}$, the model is encouraged to reconstruct x :

$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), x, c, t} [\|\epsilon - \epsilon_{\phi}(x'_t|x, \emptyset, t)\|^2], \quad (9)$$

$x'_t = \sqrt{\alpha_t}x + \sqrt{1 - \alpha_t}\epsilon$.

3. Unconditional distribution: Without conditioning input from either x_w or x_l , the model generates images drawn

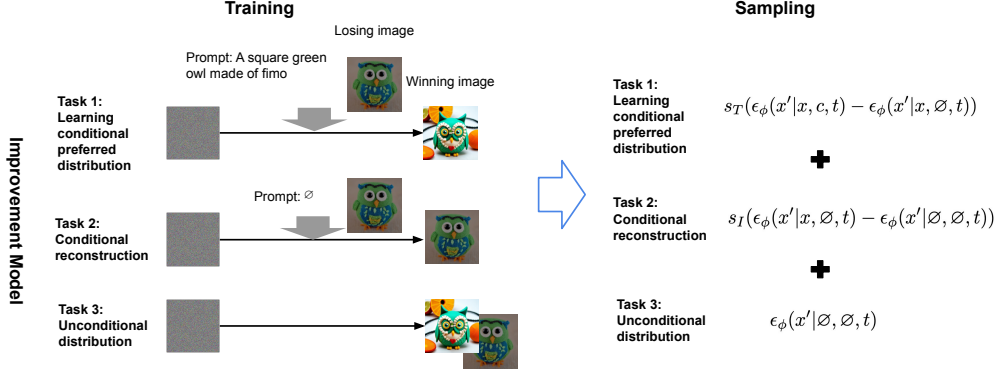


Figure 2. The overview of the training and sampling pipeline of the proposed improvement model. The left side diagram demonstrates the three tasks used for training the improvement model. The three tasks are co-trained together to make the model both learn the generic capability of improving image toward the preferred distribution represented in the offline dataset and the retain generalized image generation capability without losing diversity. The right side diagram shows the sampling strategy of the improvement model where the diffusion score of the improvement images are combined from the three tasks learned before.

from a distribution encompassing both winning and losing images:

$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), x, c, t} [\|\epsilon - \epsilon_\phi(x''_t | \emptyset, \emptyset, t)\|^2], \quad (10)$$

where $x''_t = \sqrt{\alpha_t}x + \sqrt{1 - \alpha_t}\epsilon$, $x, c \sim \mathcal{D}$.

Interpretation. The design of the improvement model is inspired by InstructPix2Pix (Brooks et al., 2023). However, their setting cannot be directly applied here because prompts in our datasets lack specific improvement/editing instructions. Furthermore, applying their objective function to our setting could cause the sampled distribution to collapse. Consider a single-task training that solely focuses on learning the conditional distribution of $p(x_w|x_l, c)$. Without an objective to force the model to utilize x_l , it might learn to ignore x_l and instead learn an unconditional distribution $p(x_w|c)$. This is likely to happen when we are fine-tuning from a pre-trained diffusion model where image condition weights are initialized to 0. This necessitates the additional reconstruction task which aims to capture the information from the image condition. Furthermore, the difference between $\epsilon_\phi(\cdot|x_l, c, t)$ and $\epsilon_\phi(\cdot|x_l, \emptyset, t)$ provides the “improvement direction” from the losing image and the prompt. Moreover, learning the unconditional score is crucial for achieving both high conditional generation accuracy and sample diversity (Ho & Salimans, 2022). We also provide ablation on the effect of the proposed reconstruction task in Section 5.5.

3.2. Sampling from the Improvement Diffusion Model

Double classifier-free guidance. For conditional sampling from the improvement model, we adapt the double classifier-free guidance technique introduced in Instruct-

Pix2Pix (Brooks et al., 2023) and design the sampling algorithm as:

$$\begin{aligned} \bar{\epsilon}_\phi(x'|x, c, t) &= \epsilon_\phi(x'|\emptyset, \emptyset, t) \\ &\quad + s_I(\epsilon_\phi(x'|x, \emptyset, t) - \epsilon_\phi(x'|\emptyset, \emptyset, t)) \\ &\quad + s_T(\epsilon_\phi(x'|x, c, t) - \epsilon_\phi(x'|x, \emptyset, t)), \end{aligned} \quad (11)$$

where x' is the output, s_T is the text guidance weight, s_I is the image guidance weight, and c is the text prompt. The first term $\epsilon_\phi(x'|\emptyset, \emptyset, t)$ is to sample without any condition as the standard diffusion model. The second term $s_I(\epsilon_\phi(x'|x, \emptyset, t) - \epsilon_\phi(x'|\emptyset, \emptyset, t))$ is to sample from the image only condition to reconstruct the input images. It helps to regularize the divergence of the output from the input images. The last term $s_T(\epsilon_\phi(x'|x, c, t) - \epsilon_\phi(x'|x, \emptyset, t))$ is to sample from both the image and text condition to improve from losing images to winning ones. The overall sampling algorithm of the improvement model is illustrated on the right side of Figure 2.

Roles of the guidance weights. To further refine the sampling process, we utilize two guidance weights: text guidance weight s_T and image guidance weight s_I . The text guidance weight s_T determines the strength of the improvement direction - a larger s_T value leads to more significant alignment with text prompt. Meanwhile, the image guidance weight s_I controls how closely the output image resembles the input condition image, i.e., increasing s_I enforces greater similarity of input and output images.

3.3. Improvement-Guided Iterative DPO for Diffusion Model Fine-tuning

The objective function. Building on the improvement diffusion model $p_{\phi}^{\dagger}(x_w|x_l, c)$, we can sample pairs of preference images x_w and x_l , where x_l are generated from the current diffusion model as the losing image and x_w output from the improvement model as the winning image. These online sampled pairs provide data for optimizing the diffusion model using the DPO objective function below:

$$\begin{aligned} & -\mathbb{E}_{c \sim \mathcal{D}, x_l \sim p_{\theta}(\cdot|c), x_w \sim p_{\phi}^{\dagger}(\cdot|x_l, c), t} [\log \sigma(-\beta(\|\epsilon^w - \epsilon_{\theta}(x_t^w, t)\|_2^2 \\ & + \|\epsilon^l - \epsilon_{\text{ref}}(x_t^l, t)\|_2^2 - \|\epsilon^w - \epsilon_{\text{ref}}(x_t^w, t)\|_2^2 \\ & - \|\epsilon^l - \epsilon_{\theta}(x_t^l, t)\|_2^2))], \end{aligned} \quad (12)$$

where the losing images x_l are from the current output of the model, and the winning images x_w are from the improvement model conditioning on x_l and c , constructing a new preference dataset. After one iteration of optimization, we can regenerate new preference data from the current diffusion model and the improvement model, and perform DPO training iteratively (details in Algorithm 1). An illustration of the iterative training pipeline is in Figure 1.

Comparison with SPIN. Here we compare our method with SPIN, a self-play method that can be applied to diffusion models (Yuan et al., 2024a). The iterative objective function of SPIN aims to move the model’s output distribution closer to a target distribution. However, a key limitation of this approach is its strong reliance on the quality of the preferred responses. SPIN uses these preferred responses, along with the prompt set, to construct an SFT dataset, while discarding the losing responses. This strategy assumes that the preferred distribution is near-optimal. If this assumption doesn’t hold, the model risks falling into a suboptimal area. Furthermore, the output distribution is still constrained by the available preferred responses in the training dataset, potentially limiting the outputs’ diversity. In contrast, we learn the improvement direction from the preference dataset while retaining information from the losing distribution. By iteratively applying this learned improvement direction, we can optimize the model towards better performances. Thus, our method could surpass SPIN models, achieving both higher alignment and higher diversity.

4. Related Work

Variants of DPO. Direct preference optimization (DPO) (Rafailov et al., 2024) is developed to optimize the generation policy with the offline preference dataset. It eliminates the dependency on the explicit reward model. However, the optimal solution derived from the Bradley-Terry (BT) model makes DPO prone to weakening the regularization and overfitting to the offline training dataset. Azar et al.

Algorithm 1 Improvement-guided iterative DPO training.

Input: Improvement model p_{ϕ}^{\dagger} , prompt set \mathcal{D}_c , model p_{θ} , number of iterations T_{iter} , number of samples n , training batch size b , text guidance weight s_T , image guidance weight s_I , steps per iteration T_{train}

for $t_{\text{iter}} \in [1, T_{\text{iter}}]$ **do**

Randomly sample n images from p_{θ} conditioned on \mathcal{D}_c , and construct \mathcal{D}_l

Randomly sample n images from p_{ϕ}^{\dagger} conditioned on \mathcal{D}_c and \mathcal{D}_l . With guidance weights s_T and s_I , construct \mathcal{D}_w

for $t_{\text{train}} \in [1, T_{\text{train}}]$ **do**

Compute an estimation of gradient using Equation (12) with batch size b , and update θ

end for

end for

Output: Fine-tuned model p_{θ}

(2024) propose the IPO by introducing the identity function into the generic Ψ PO framework and derive an efficient optimization process and achieve improved performance than DPO. Meng et al. (2024) argue for the effectiveness of the reference model regularization in DPO. They therefore propose the simple preference optimization (SimPO) method that bypasses the reference model regularization and introduces a reward margin to the optimization objective to better approximate the noisy preference dataset. Their approach also shows improved performance over DPO. Hong et al. (2024) also argue about impediments in optimizing the reference model under distributional discrepancy and propose the margin-aware preference optimization (MaPO) method to replace KL regularization on the reference model with an amplification factor defined by the trained policy’s likelihood estimation. These DPO variants explore the challenges of distribution discrepancy between the reference model and model under optimization. They optimize with the offline samples which is verified to be less efficient than on-policy sampling (Tajwar et al., 2024).

Iterative DPO and self-play methods. To understand and address the limitations of offline training associated with DPO, recent works have investigated the performance gap between online and offline training methods (Tajwar et al., 2024; Tang et al., 2024). Their findings indicate that online training can lead to better generation, and is beneficial when high-reward responses have a low likelihood under the pretrained model. Accordingly, several works have proposed iterative DPO methods that train DPO using online samples generated by the improved policy (Guo et al., 2024; Xu et al., 2023a; Xiong et al., 2023). However, they require a reward model to label the online samples. To eliminate the dependence on reward models, researchers have developed

self-play or self-improvement methods. For example, Yuan et al. (2024b) use the language model itself to provide the reward signal, and Chen et al. (2024) treat self-generated responses as losing to human demonstrations for iterative improvement. More recently, Choi et al. (2024); Wu et al. (2024) reformulate these ideas under the constant-sum two-player game framework (Munos et al., 2023; Swamy et al., 2024), and propose algorithms to find the approximate Nash equilibrium. Our work proposes a self-improvement method for text-to-image diffusion models, which has been under-explored.

Aligning diffusion models with human preferences. Inspired by the success of RLHF and DPO in fine-tuning language models, recent works have explored applications in aligning diffusion models with human preferences. RLHF-based methods maximize a reward score given by a separately trained reward model (Radford et al., 2021; Lee et al., 2023; Xu et al., 2023b; Wu et al., 2023; Kirstain et al., 2023). For differentiable rewards, reward maximization can be done by backpropagating the reward function gradient through the denoising process (Clark et al., 2024; Prabhudesai et al., 2023). For black-box reward functions, DDPO (Black et al., 2024) and DPOK (Fan et al., 2023) propose PPO-based RL fine-tuning. PRDP (Deng et al., 2024) further improves training stability on large-scale datasets by converting reward maximization to an equivalent reward difference prediction objective. However, RLHF-based methods generally have a complicated pipeline involving reward model training, and are prone to reward hacking. These issues can be partially mitigated by DPO-based methods, such as Diffusion-DPO (Wallace et al., 2024) and SPIN-Diffusion (Yuan et al., 2024a), which directly fine-tune the diffusion model from offline preference datasets without requiring reward models. However, their performance can be limited due to a lack of online training. Our approach combines the benefits of online training from RLHF and the simplicity from DPO.

5. Experiments

5.1. Experimental Setup

5.1.1. TRAINING

Model and Dataset. We use the Pick-a-pic (Kirstain et al., 2023) training dataset as the offline preference dataset, following Diffusion-DPO (Wallace et al., 2024). For the improvement model, we add 4 channels to the first convolutional layer of the UNet, and initialize the weights from Stable Diffusion 1.5 (Rombach et al., 2022) following (Brooks et al., 2023). For iterative DPO training, we fine-tune the model initialized from the third iteration in SPIN (Yuan et al., 2024a).

Hyperparameters. For the improvement model training, we use AdamW with a learning rate 10^{-4} , and train up to 200K steps with batch size 2048, and sample from it with $s_T = 3.5$, $s_I = 3.0$. For iterative DPO training, we train for 3 iterations, and for each iteration, we first generate 38400 pairs of preference data, and train for 5k steps for each iteration with the batch size 2048 and learning rate 10^{-4} , $\beta = 2000$, with SD 1.5 as the reference model.

5.1.2. EVALUATION

Prompt sets. We use two prompt sets for evaluation: We randomly sample 500 unique prompts from the training set to reflect the model’s performance on the training set. We also use the 500 unique prompts from the test dataset in Pick-a-pic v2 as the test set. We sample 64 random images from each prompt.

Metrics. We evaluate our method against DPO and other baselines using a comprehensive set of metrics. For quality assessment, we employ PickScore (Kirstain et al., 2023), Human Preference Score v2 (HPSv2) (Wu et al., 2023) and Aesthetic score (Schuhmann et al., 2022), which capture different aspects of image quality and alignment with human preferences. To ensure that our method not only improves quality but also maintains diversity in generated images, we use the Vendi score (Friedman & Dieng, 2023) to measure the diversity.

Baseline methods. We compare our method with the base model SD 1.5, and DPO-based methods: Diffusion-DPO and SPIN. Notice that Diffusion-DPO, SPIN, and our method all share the same data assumption: using the offline preference dataset only without the need for extra annotations or feedback from reward models.

5.2. Reward Evaluation

We report the results of Pickscore, HPSv2, and Aesthetic score in Table 1. Our iterative training can further improve the SPIN model with prompts in both training and test sets in terms of the surrogate metrics of human preferences. It further proves that our iterative training with an improvement model can surpass the upper bound in SPIN training. We also visualize samples from the fine-tuned model in Figure 6, where our model output can generate the samples more aligned with the prompt than the SPIN model.

5.3. Human Evaluation

We conduct human evaluation using the prompts in the Pick-a-pic test set with two metrics: visual quality of the image and text-image alignment (see details in Appendix B). We present the results in Figure 5 where our fine-tuned model consistently outperforms the baseline. We also present visu-



Figure 3. Visualization of sampled images from the baseline and fine-tuned diffusion models, where our fine-tuning improves the text-image alignment. Prompts (left to right): 1. A guinea pig riding a motorcycle; 2. Of the lunar module landing on a hydrogen lake on Titan, through a foggy yellow smog; 3. Earth from Mars; 4. Illustration cartoon of a leprechaun gnome with a rainbow hat stuck at the bottom of a rock pit.

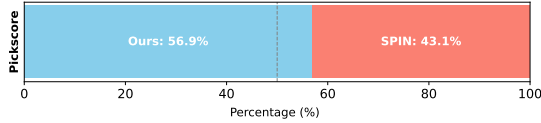


Figure 4. Automatic win-rate of our model against SPIN calculated by Pickscore, evaluated on Pick-a-pic test set.

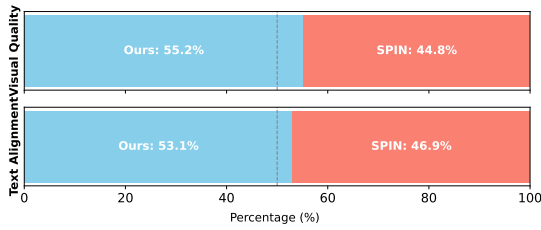


Figure 5. Win-rate evaluated on visual quality and text alignment respectively conducted by human.

alizations in Figure 3 for text-image alignment and Figure 6 for visual quality.

5.4. Effect of the Online Samples

Here we present the effect of the number of online samples used for iterative training. From Table 2, we find that more online samples can lead to higher Pickscore from prompts in both training and test sets. This verifies that the key to successful iterative training is the online samples generated from our improvement model: more online samples would

Table 1. Evaluation of Pickscore, HPSv2, Aesthetic score and Vendi score. Our model achieves improved reward scores without sacrificing diversity compared with SPIN and Diffusion-DPO.

Score	Method	Training subset	Test set
Pickscore (\uparrow)	SD 1.5	20.46	20.74
	Diffusion-DPO	20.80	21.05
	SPIN	<u>21.15</u>	<u>21.41</u>
	Iterative (Ours)	21.22	21.47
HPSv2 (\uparrow)	SD 1.5	26.65	26.90
	Diffusion-DPO	26.96	27.19
	SPIN	27.39	27.57
	Iterative (Ours)	27.44	27.61
Aesthetic (\uparrow)	SD 1.5	5.48	5.42
	Diffusion-DPO	5.55	5.49
	SPIN	<u>5.92</u>	<u>5.86</u>
	Iterative (Ours)	5.95	5.88
Vendi score (\uparrow)	SD 1.5	2.61	2.64
	Diffusion-DPO	2.44	2.47
	SPIN	2.43	2.48
	Iterative (Ours)	<u>2.45</u>	<u>2.49</u>

lead to better results.

5.5. Evaluation of the Improvement Model

In this section, we provide the evaluation of the improvement model, by using the SD 1.5 baseline to generate the losing images as the image condition for the improvement model. We find that the improvement model can achieve significant improvements on the training set, but the gen-



Figure 6. Visualization of the sampled images from the baseline and the fine-tuned diffusion models, where our fine-tuning improves the visual quality. Prompts (left to right): 1. A cyborg on the ocean; 2. Cute grey cat, digital oil painting by Monet; 3. Gray French bulldog; 4. hummingbird.

Table 2. Evaluation on the effect of the online samples. The values presented are Pickscore.

Number of online samples	Training set	Test set
2560	21.10	21.27
12800	21.13	21.36
38400	21.22	21.47

Table 3. Ablation study: Evaluation of the improvement model with or without reconstruction training. The values presented are Pickscore.

Method	Training set	Test set
With reconstruction	21.24	21.38
Without reconstruction	21.17	21.29

eralization ability on the test set is worse than the iterative model, which implies why we do not consider using it as an inference-time model. The improvement model modifies the original architecture of SD 1.5 and is trained with different tasks than text-to-image generation. Thus the generalization ability on test prompts may not be as good as fine-tuned diffusion model. In the iterative training, we reuse the same training prompts and do not use the improvement model on unseen prompts. The iteratively trained model with the improvement model can therefore achieve better generalization ability on the test set. We also include samples generated by the improvement model in Appendix C.

6. Discussion and Limitation

Note that the gap between SPIN and our method depends on the specific structure of the preferences dataset. If all winning images in the preference set are near-optimal, there is little space for improvement with our improvement model and iterative training. However, if the winning images contain a diverse range from sub-optimal to optimal, SPIN can only get mediocre quality at best. In contrast, our method that learns the improvement direction can outperform SPIN. Due to a lack of resources, we use the open-source benchmark dataset instead of creating a more diverse dataset for losing images that could potentially lead to larger improvements from SPIN.

References

- Azar, M. G., Guo, Z. D., Piot, B., Munos, R., Rowland, M., Valko, M., and Calandriello, D. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR, 2024.
- Black, K., Janner, M., Du, Y., Kostrikov, I., and Levine, S. Training diffusion models with reinforcement learning. In *International Conference on Learning Representations*, 2024.
- Brooks, T., Holynski, A., and Efros, A. A. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18392–18402, 2023.
- Chen, Z., Deng, Y., Yuan, H., Ji, K., and Gu, Q. Self-play

- fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024.
- Choi, E., Ahmadian, A., Geist, M., Pietquin, O., and Azar, M. G. Self-improving robust preference optimization. *arXiv preprint arXiv:2406.01660*, 2024.
- Clark, K., Vicol, P., Swersky, K., and Fleet, D. J. Directly fine-tuning diffusion models on differentiable rewards. In *International Conference on Learning Representations*, 2024.
- Deng, F., Wang, Q., Wei, W., Hou, T., and Grundmann, M. PRDP: Proximal reward difference prediction for large-scale reward finetuning of diffusion models. In *CVPR*, 2024.
- Fan, Y., Watkins, O., Du, Y., Liu, H., Ryu, M., Boutilier, C., Abbeel, P., Ghavamzadeh, M., Lee, K., and Lee, K. DPOK: Reinforcement learning for fine-tuning text-to-image diffusion models. In *Advances in Neural Information Processing Systems*, 2023.
- Friedman, D. and Dieng, A. B. The vendi score: A diversity evaluation metric for machine learning. *Transactions on machine learning research*, 2023.
- Guo, S., Zhang, B., Liu, T., Liu, T., Khalman, M., Llinares, F., Rame, A., Mesnard, T., Zhao, Y., Piot, B., et al. Direct language model alignment from online AI feedback. *arXiv preprint arXiv:2402.04792*, 2024.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Hong, J., Paul, S., Lee, N., Rasul, K., Thorne, J., and Jeong, J. Margin-aware preference optimization for aligning diffusion models without reference. *arXiv preprint arXiv:2406.06424*, 2024.
- Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., and Levy, O. Pick-a-Pic: An open dataset of user preferences for text-to-image generation. In *Advances in Neural Information Processing Systems*, 2023.
- Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C., Abbeel, P., Ghavamzadeh, M., and Gu, S. S. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- Meng, Y., Xia, M., and Chen, D. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*, 2024.
- Munos, R., Valko, M., Calandriello, D., Azar, M. G., Rowland, M., Guo, Z. D., Tang, Y., Geist, M., Mesnard, T., Michi, A., et al. Nash learning from human feedback. *arXiv preprint arXiv:2312.00886*, 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.
- Prabhudesai, M., Goyal, A., Pathak, D., and Fragkiadaki, K. Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021.
- 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. *Advances in Neural Information Processing Systems*, 36, 2024.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Rosset, C., Cheng, C.-A., Mitra, A., Santacroce, M., Awadallah, A., and Xie, T. Direct nash optimization: Teaching language models to self-improve with general preferences. *arXiv preprint arXiv:2404.03715*, 2024.
- Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., 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.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Swamy, G., Dann, C., Kidambi, R., Wu, Z. S., and Agarwal, A. A minimaximalist approach to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056*, 2024.
- Tajwar, F., Singh, A., Sharma, A., Rafailov, R., Schneider, J., Xie, T., Ermon, S., Finn, C., and Kumar, A. Preference fine-tuning of llms should leverage suboptimal, on-policy data. In *Forty-first International Conference on Machine Learning*, 2024.

- Tang, Y., Guo, D. Z., Zheng, Z., Calandriello, D., Cao, Y., Tarassov, E., Munos, R., Pires, B. Á., Valko, M., Cheng, Y., et al. Understanding the performance gap between online and offline alignment algorithms. *arXiv preprint arXiv:2405.08448*, 2024.
- Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A., Pushwalkam, S., Ermon, S., Xiong, C., Joty, S., and Naik, N. Diffusion model alignment using direct preference optimization. In *CVPR*, 2024.
- Wu, X., Hao, Y., Sun, K., Chen, Y., Zhu, F., Zhao, R., and Li, H. Human Preference Score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- Wu, Y., Sun, Z., Yuan, H., Ji, K., Yang, Y., and Gu, Q. Self-play preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*, 2024.
- Xiong, W., Dong, H., Ye, C., Zhong, H., Jiang, N., and Zhang, T. Gibbs sampling from human feedback: A provable KL-constrained framework for RLHF. *arXiv preprint arXiv:2312.11456*, 2023.
- Xiong, W., Dong, H., Ye, C., Wang, Z., Zhong, H., Ji, H., Jiang, N., and Zhang, T. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In *Forty-first International Conference on Machine Learning*, 2024.
- Xu, J., Lee, A., Sukhbaatar, S., and Weston, J. Some things are more cringe than others: Preference optimization with the pairwise cringe loss. *arXiv preprint arXiv:2312.16682*, 2023a.
- Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., and Dong, Y. ImageReward: Learning and evaluating human preferences for text-to-image generation. In *Advances in Neural Information Processing Systems*, 2023b.
- Xu, S., Fu, W., Gao, J., Ye, W., Liu, W., Mei, Z., Wang, G., Yu, C., and Wu, Y. Is dpo superior to ppo for llm alignment? a comprehensive study. In *Forty-first International Conference on Machine Learning*, 2024.
- Yuan, H., Chen, Z., Ji, K., and Gu, Q. Self-play fine-tuning of diffusion models for text-to-image generation. *arXiv preprint arXiv:2402.10210*, 2024a.
- Yuan, W., Pang, R. Y., Cho, K., Li, X., Sukhbaatar, S., Xu, J., and Weston, J. E. Self-rewarding language models. In *Forty-first International Conference on Machine Learning*, 2024b.
- Zhang, Y., Tzeng, E., Du, Y., and Kislyuk, D. Large-scale reinforcement learning for diffusion models. *arXiv preprint arXiv:2401.12244*, 4, 2024.

A. More Visualizations

Here we provide more image samples from both our trained model and SPIN in Figures 3, 6, 7, 8 and 9, where the outputs from our model are generally more aligned with the text description and have better image quality than SPIN.



Figure 7. Visualization of the sampled images from the baseline and the fine-tuned diffusion models. Prompts (left to right): 1. Purple cat eating cake; 2. A white cat wearing a red hat holding sticks; 3. Cute simple rabbit lineart; 4. A diamond ring on a girls hand; 5. Tattoo ideas for an introvert. The outputs from our model have better **visual quality** in general.

B. Human Evaluation Details

Here we provide details in our human evaluation in Section 5.3. The total number of prompts in Pick-a-pic test set is 500, and we filtered out 43 prompts that are not suitable for work. We have 47 human raters to label the results from 457 unique prompts (for each prompt there is one side-by-side image pair generated from the same seed for both models we compare). The total number of ratings for image quality is 4050, where each rater gives 86.2 ratings on average, and each unique prompt repeats 8.86 times on average. The total number of ratings for text-image alignment is 3009, where each rater gives 64.0 ratings on average, and each unique prompt repeats 6.58 times on average.

C. Samples from the Improvement Model

We also provide visualizations of the output of the improvement model in Figure 10, starting from an image condition generated from the based SD 1.5 model and then repeatedly apply the model to generate new images. The prompts are from the Pick-a-pic test set. We can observe that the output of the improvement model from the first iteration is significantly improved from the initial image, while the quality of output from the second iteration is slightly better or similar comparing to the first iteration. Similar to the discussion in Section 6, if the winning images in the training set are mostly near-optimal ones, the room for improvement after applying the first iteration might be limited. Nevertheless, such improvement model could still effectively guide the iterative DPO process. In our paper, we also show the improvement from SPIN is possible given such training datasets. Exploring the potential benefit of our method when training on more diverse datasets could be an interesting future work.

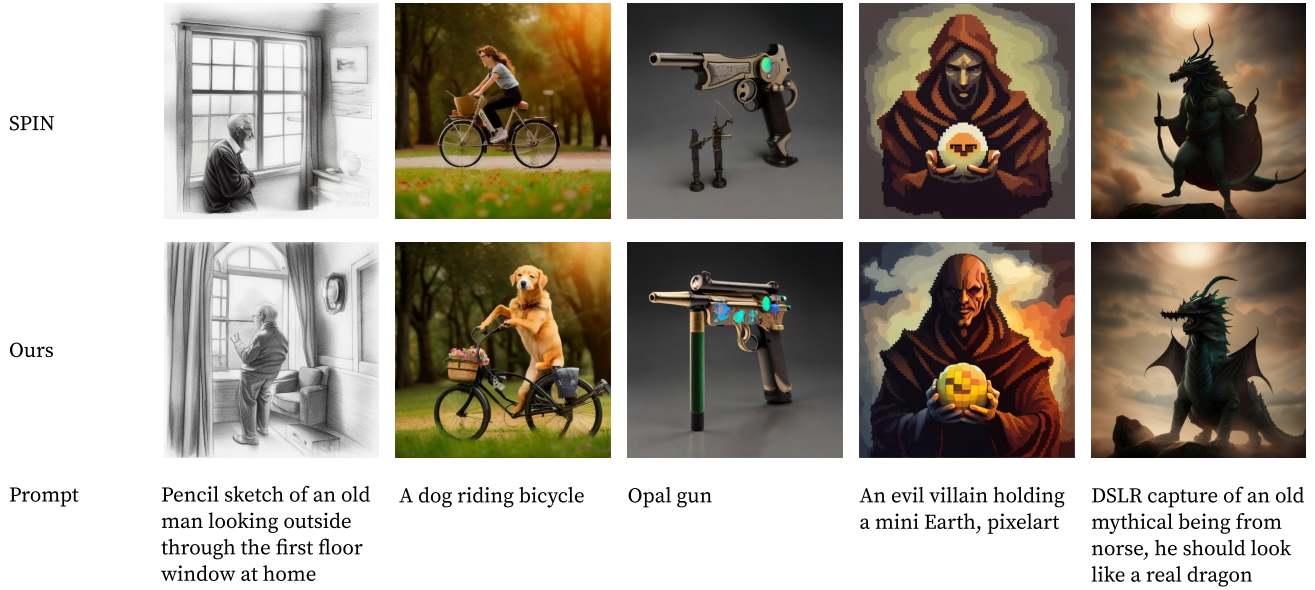


Figure 8. Visualization of the sampled images from the baseline and the fine-tuned diffusion models. Prompts (left to right): 1. Pencil sketch of an old man looking outside through the first floor window at home; 2. A dog riding bicycle; 3. Opal gun; 4. An evil villain holding a mini Earth, pixelart; 5. DSLR capture of an old mythical being from norse, he should look like a real dragon. The outputs from our model have better **text-image alignment** in general.

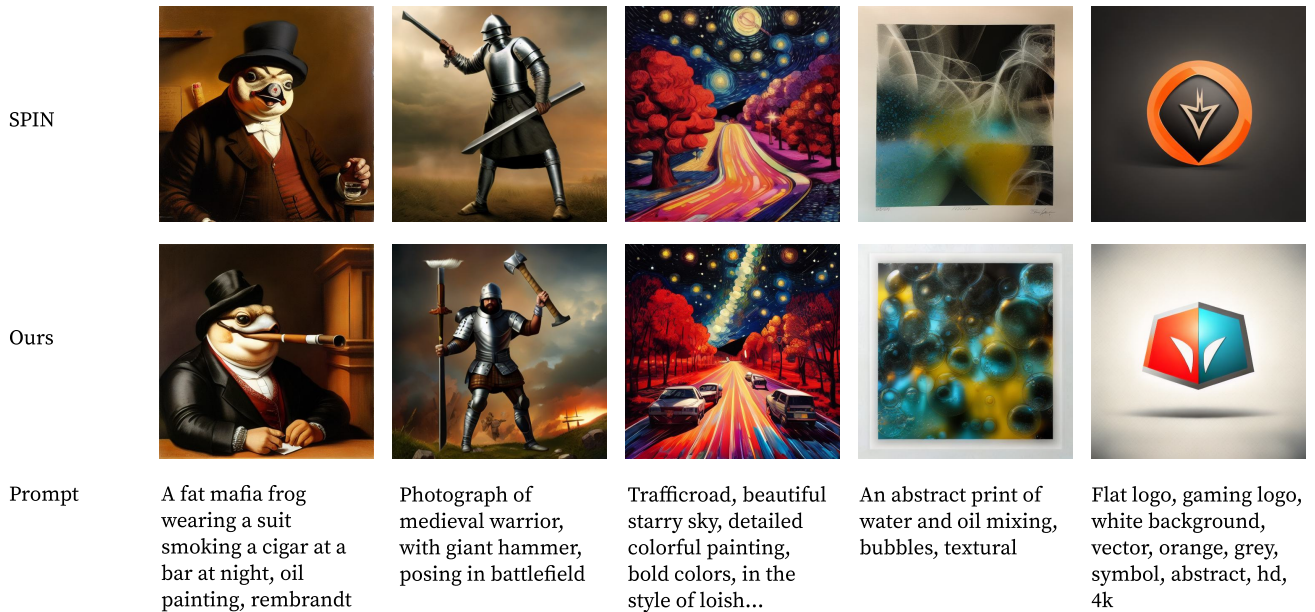


Figure 9. Visualization of the sampled images from the baseline and the fine-tuned diffusion models. Prompts (left to right): 1. A fat mafia frog wearing a suit smoking a cigar at a bar at night, oil painting, rembrandt; 2. Photograph of medieval warrior, with giant hammer, posing in battlefield; 3. Trafficroad, beautiful starry sky, detailed colorful painting, bold colors, in the style of loish, artist Antonio Ligabue, artist Djanira, by van gogh, artist Marija Prymatschenko and william morris print, sharp lines, intricate, fine black outlines, light and shadow, octane render, Ultra-violet ink painting, by Dan Mumford, inspired by cyberpunk and retro-futuristic elements, impressionism, featuring neon blue and purple hues, bold lines, and dynamic composition, expressionism, created with Copic markers and black light ink, evokes a sense of mystery, sophistication, and otherworldly beauty; 4. An abstract print of water and oil mixing, bubbles, textural; 5. Flat logo, gaming logo, white background, vector, orange, grey, symbol, abstract, hd, 4k. The outputs from our model have better **text-image alignment** in general.

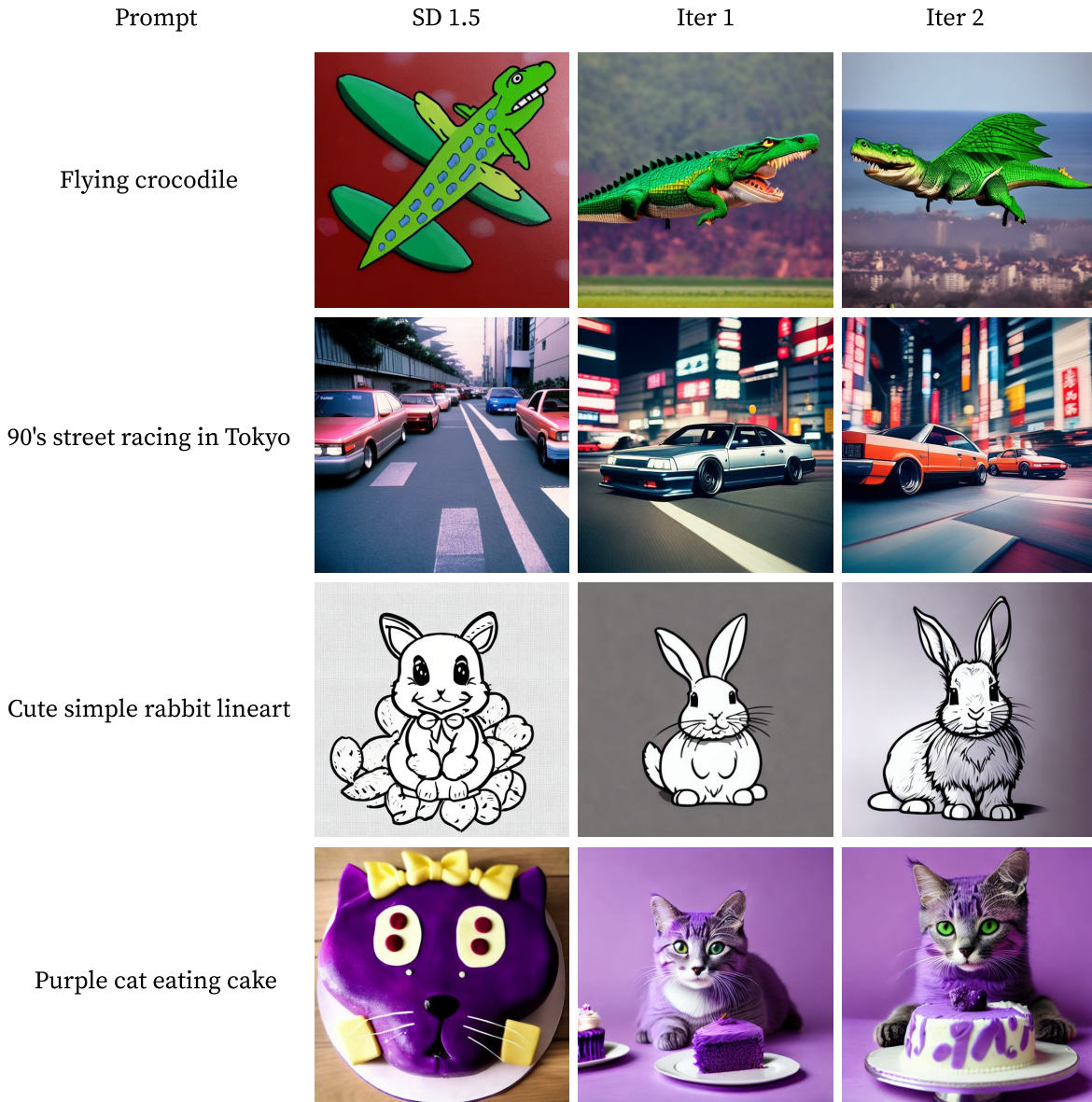


Figure 10. Outputs from the improvement model. Prompts (top to bottom): 1. Flying crocodile; 2. 90's street racing in Tokyo; 3. Cute simple rabbit lineart; 4. Purple cat eating cake.