

---

# GUIDING THE TRAINING OF DIFFUSION MODELS: A REPLICATION OF DDPO

---

**Gao, Qiaosong**  
New York University, Shanghai  
qg2089@nyu.edu

**Shi, Xingyu**  
New York University, Shanghai  
xs2337@nyu.edu

**Zhang, Haoqi**  
New York University, Shanghai  
hz3223@nyu.edu

## ABSTRACT

Diffusion models have emerged as a class of effective generative models, trained with an approximation to the log-likelihood objective. However, this objective sometimes doesn't best reflect the objective of empirical use cases of diffusion models such as human-perceived image quality. The paper we want to replicate: *Training Diffusion Models with Reinforcement Learning* [1], aim to improve the capabilities of diffusion models for such empirical objectives through reinforcement learning methods. They describe denoising as a multi-step decisionmaking problem, which enables a class of policy gradient algorithms, which they refer to as denoising diffusion policy optimization (DDPO). We replicate DDPO, and based on this, propose our own insights. The code for our replication and implementation can be found at [Fall-2025-Reinforcement-Learning-Final-Project](#).

## 1 Introduction

Diffusion models are widely used to generate high-quality images and videos. A key challenge in training empirically effective diffusion models is that the evaluation metrics for generated samples (such as perceptual image quality) are not directly defined in terms of the likelihood-based objectives that diffusion models optimize during training.

To better train diffusion models that align with the objectives of empirical tasks, Black et al. [1] frames denoising as a multi-step decision-making task, introducing reinforcement learning methods. Specifically, they present a policy gradient algorithm, or Denoising Diffusion policy Optimization (DDPO), which is able to optimize a diffusion model for downstream tasks using different reward functions. Various reward functions are also designed for different downstream tasks. The downstream tasks they aim to optimize are: Compressibility, Aesthetic Quality, and Prompt-Image Alignment. DDPO and these reward functions are used for training on different downstream objectives and demonstrate their effectiveness compared to alternative reward-weighted likelihood methods.

We've successfully replicated the DDPO algorithm, and achieved similar results on similar downstream objectives. We've also proposed our improvement based on one class of the DDPO algorithm and achieved faster convergence and avoided overoptimization in some tasks.

## 2 Environment

The environment of the DDPO reinforcement learning method is text-to-image diffusion. Text-to-image diffusion "serves as a valuable test environment for reinforcement learning due to the availability of large pretrained models and the versatility of using diverse and visually interesting reward functions [1]. "

In this section, we describe the environment of this reinforcement learning method, including:

1. The Markov Decision Proces (MDP) formulation
2. Description of a selection of downstream tasks and corresponding reward functions

## 2.1 Markov Decision Process Formulation

A Markov decision process (MDP) is a formalization of sequential decision-making problems. An MDP is defined by a tuple  $(S, A, \rho_0, P, R)$ , in which  $S$  is the state space,  $A$  is the action space,  $\rho_0$  is the distribution of initial states,  $P$  is the transition kernel, and  $R$  is the reward function. At each timestep  $t$ , the agent observes a state  $s_t \in S$ , takes an action  $a_t \in A$ , receives a reward  $R(s_t, a_t)$ , and transitions to a new state  $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$ . An agent acts according to a policy  $\pi(a_t | s_t)$ .

As the agent acts in the MDP, it produces trajectories, which are sequences of states, actions  $\tau = (s_0, a_0, s_1, a_1, \dots, s_T, a_T)$ . The reinforcement learning (RL) objective is for the agent to maximize  $J_{\text{RL}}(\pi)$ , the expected cumulative reward over trajectories sampled from its policy:

$$J_{\text{RL}}(\pi) = \mathbb{E}_{\tau \sim p(\tau | \pi)} \left[ \sum_{t=0}^T R(s_t, a_t) \right].$$

## 2.2 MDP formulation of the denoising process of diffusion models

Now, we consider the **MDP formulation of the denoising process of diffusion models**. We cast conditional diffusion sampling as an episodic finite-horizon Markov decision process (MDP). Each episode corresponds to one reverse diffusion trajectory that starts from pure noise and ends at a synthesized image.

**State space.** Let  $c \in \mathcal{C}$  denote a text prompt and  $x_t \in \mathbb{R}^d$  the latent (noisy image) at diffusion step  $t \in \{0, \dots, T\}$ . We define the state at time  $t$  as

$$s_t \triangleq (c, t, x_t) \in \mathcal{S} = \mathcal{C} \times \{0, \dots, T\} \times \mathbb{R}^d.$$

The agent therefore observes the current prompt, the current diffusion time index, and the current noisy latent.

**Action space.** At each step  $t > 0$  the agent chooses the next latent  $x_{t-1}$ . The action space is the same Euclidean space as the latent:

$$a_t \triangleq x_{t-1} \in \mathcal{A} = \mathbb{R}^d.$$

**Policy.** The policy is given by the reverse diffusion model  $p_\theta$ :

$$\pi_\theta(a_t | s_t) = \pi_\theta(x_{t-1} | c, t, x_t) \triangleq p_\theta(x_{t-1} | x_t, c, t),$$

which is typically a Gaussian with mean  $\mu_\theta(x_t, c, t)$  and fixed variance  $\sigma_t^2 I$  (possibly after classifier-free guidance).

**Transition dynamics.** The environment transition simply decreases the time index and copies the chosen latent:

$$P(s_{t-1} | s_t, a_t) = P((c', t', x'_{t-1}) | (c, t, x_t), x_{t-1}) = \delta_c(c') \delta_{t-1}(t') \delta_{x_{t-1}}(x'_{t-1}),$$

i.e., the prompt  $c$  is fixed throughout the episode, the time index counts down from  $T$  to 0, and the next latent is deterministically exactly the chosen action.

**Initial state distribution.** At the beginning of each episode we sample a prompt and an initial noise:

$$c \sim p(c), \quad x_T \sim \mathcal{N}(0, I),$$

and set the initial state to

$$s_T = (c, T, x_T), \quad \rho_0(s_T) = p(c) \mathcal{N}(x_T; 0, I).$$

**Reward.** The episode terminates after producing  $x_0$  at  $t = 0$ . We only assign a terminal reward based on the final image and prompt:

$$R(s_t, a_t) = \begin{cases} r(x_0, c), & t = 0, \\ 0, & t > 0, \end{cases}$$

where  $r(x_0, c)$  is given by an external reward model (e.g., aesthetic score, compression ratio, or a VLM).

**Objective.** Let  $\tau = (s_T, a_T, \dots, s_0, a_0)$  denote a full trajectory induced by  $\pi_\theta$ . The DDRL objective is to maximize the expected terminal reward

$$J_{\text{DDRL}}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r(x_0, c)].$$

**One-step RWR view.** For reward-weighted regression (RWR), this MDP can be collapsed to a one-step formulation:

$$s \triangleq c, \quad a \triangleq x_0, \quad \pi_\theta(a | s) \triangleq p_\theta(x_0 | c), \quad R(s, a) \triangleq r(x_0, c),$$

which leads to a weighted minimization of the diffusion MSE loss objective for  $p_\theta(x_0 | c)$ .

## 2.3 A Selection of Tasks and Corresponding Reward Functions

This paper mainly evaluate the methods on the following selection of downstream of tasks and corresponding reward functions:

1. Compressibility and Incompressibility
2. Aesthetic Quality

### 2.3.1 Compressibility and Incompressibility

Images are rarely captioned with their file size, making it impossible to specify a desired file size via prompting. This limitation makes reward functions based on file size a convenient case study: they are simple to compute, but not controllable through the conventional methods of likelihood maximization and prompt engineering. We fix the resolution of diffusion model samples at 512x512, such that the file size is determined solely by the compressibility of the image. We define two tasks based on file size: compressibility, in which the file size of the image after JPEG compression is minimized, and incompressibility, in which the same measure is maximized.

We adopt this task setting and reward function (image file size) in our replication.

### 2.3.2 Aesthetic Quality

To capture a reward function that would be useful to a human user, we define a task based on perceived aesthetic quality. We use the LAION aesthetics predictor [2], which is trained on 176,000 human image ratings. Annotations range between 1 and 10, with the highest-rated images mostly containing artwork.

We adopt this task setting and reward function (LAION aesthetics predictor rating) in our replication.

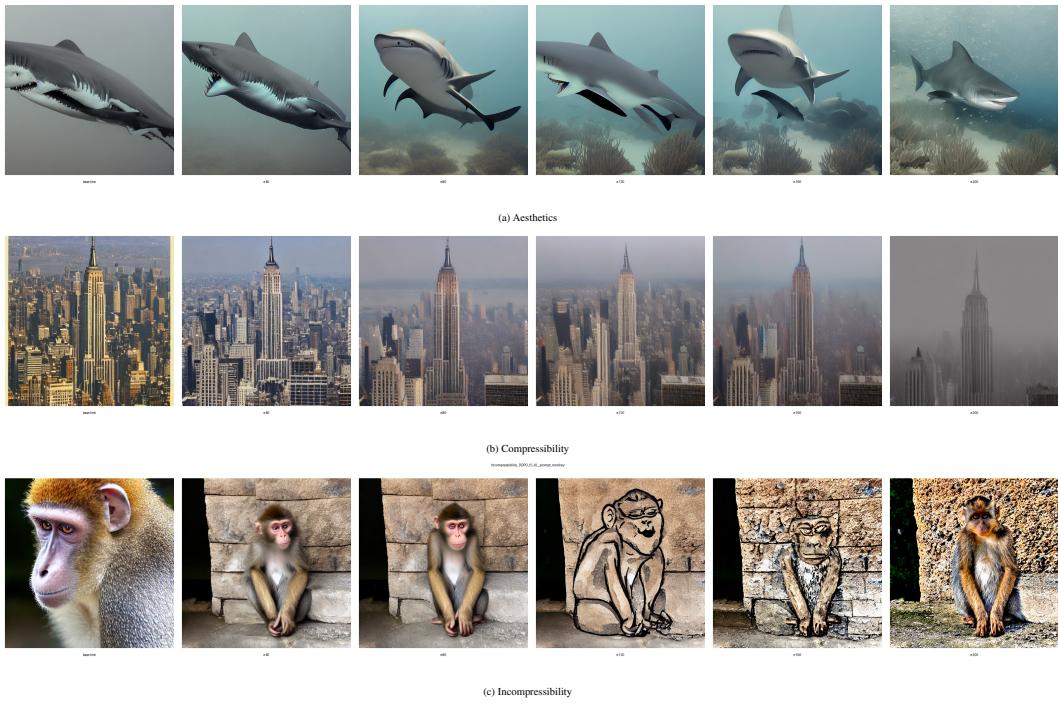


Figure 1: Training progression for different tasks



Figure 2: Training progression for compressibility with DDPO IS and SF

### 3 Algorithm and Proposed Improvement

This section includes the algorithms we choose to replicate and run. In Black et al.'s paper, they proposed the following four algorithms[1]:

1. Reward-Weighted Regression (RWR) (See section 3.1)
2. Sparse Reward-Weighted Regression (See section 3.2)
3. DDPO<sub>SF</sub> (See section 3.3)
4. DDPO<sub>IS</sub> (See section 3.4)

We also propose our improvement based on DDPO<sub>IS</sub>, introducing KL-Divergence regularization between Model in training  $p_\theta$  and a pretrained stable diffusion model  $p_{ref}$ . Through this improvement, we achieved faster convergence and avoided overoptimization in some tasks. Details and pseudo-code can be found in

5. DDPO<sub>ISKL</sub> (See section 3.5)

#### 3.1 Reward-Weighted Regression (RWR)

We have The DDRL objective:

$$J_{\text{DDRL}}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [ r(x_0, c) ].$$

Now, with training data  $x_0 \sim p_0(x_0|c)$  and an added weighting that depends on the reward  $r(x_0, c)$ . This approach can be performed for multiple rounds of alternating sampling and training, leading to an online RL method. We refer to this general class of algorithms as reward-weighted regression (RWR).

We have the following standard weighting scheme uses exponentiated rewards to ensure non-negativity,

$$w_{\text{RWR}}(x_0, c) = \frac{1}{Z} \exp (\beta r(x_0, c)),$$

where  $\beta$  is an inverse temperature and  $Z$  is a normalization constant.

We have the following pseudo-code for RWR (a policy gradient method):

---

**Algorithm 1** Denoising Diffusion with Reward-Weighted Regression (RWR)

---

```

1: Input: initial diffusion parameters  $\theta$ ; reward function  $r(x, c)$ ; prompt distribution  $p(c)$ ; timestep distribution  $p(t)$ .
2: Hyperparameters: reward temperature  $\beta$ ; number of gradient updates per iteration  $N_{\text{updates}}$ ; number of diffusion steps  $T$ ; CFG weight  $W_c$ ; number of samples per iteration  $N_{\text{samples}}$ ; number of inference steps  $N_{\text{inf}}$ ; batch size  $B$ ; noise schedule endpoints  $\beta_{\text{start}}, \beta_{\text{end}}$ ; learning rate  $\eta$ ; Adam parameters  $(\beta_1, \beta_2)$ ; Adam numerical constant  $\varepsilon$ ; weight decay  $\lambda$ ; gradient clipping norm  $C_{\text{clip}}$ .
3: repeat
4:   Initialize empty dataset  $\mathcal{D}$ 
5:   Use a diffusion scheduler with  $(T, \beta_{\text{start}}, \beta_{\text{end}})$  to generate the cumulative factors  $\{\bar{\alpha}_t\}_{t=0}^{T-1}$ 
6:   for  $i = 1$  to  $N_{\text{samples}}$  do
7:     Sample prompt  $c_i \sim p(c)$ 
8:     Run policy  $\pi_\theta(\cdot | c_i, w)$  to get a trajectory with  $N_{\text{inf}}$  time steps, the reward in each step is  $r_{\text{raw}}^{(i)} \leftarrow r(x_0^{(i)}, c_i)$ 
9:     Collect the terminal  $(c_i, x_0^{(i)}, r_{\text{raw}}^{(i)})$  to  $\mathcal{D}$ 
10:    end for
11:    for all prompts  $c$  in  $\mathcal{D}$  do
12:      Let  $\mathcal{I}_c$  be indices with prompt  $c$ 
13:       $Z_c \leftarrow \sum_{j \in \mathcal{I}_c} \exp(\beta r_{\text{raw}}^{(j)})$ 
14:      for all  $i \in \mathcal{I}_c$  do
15:         $w_i \leftarrow \exp(\beta r_{\text{raw}}^{(i)}) / Z_c$ 
16:      end for
17:    end for
18:    for  $u = 1$  to  $N_{\text{updates}}$  do
19:      Sample a minibatch  $\mathcal{B} \subset \mathcal{D}$  of size  $B$ 
20:       $L \leftarrow 0$ 
21:      for all  $i \in \mathcal{B}$  do
22:        Sample timestep  $t_i \sim p(t)$  over  $\{0, \dots, T - 1\}$ 
23:        Sample noise  $\epsilon_i \sim \mathcal{N}(0, I)$ 
24:        Construct noised latent  $x_{t_i}^{(i)} \leftarrow \sqrt{\bar{\alpha}_{t_i}} x_0^{(i)} + \sqrt{1 - \bar{\alpha}_{t_i}} \epsilon_i$ 
25:        Compute conditional prediction  $\hat{\epsilon}_\theta^{\text{cond}} \leftarrow \epsilon_\theta(x_{t_i}^{(i)}, c_i, t_i)$ 
26:        Compute unconditional prediction  $\hat{\epsilon}_\theta^{\text{uncond}} \leftarrow \epsilon_\theta(x_{t_i}^{(i)}, \emptyset, t_i)$ 
27:        Apply classifier-free guidance:

$$\hat{\epsilon}_\theta \leftarrow \hat{\epsilon}_\theta^{\text{uncond}} + W_c (\hat{\epsilon}_\theta^{\text{cond}} - \hat{\epsilon}_\theta^{\text{uncond}})$$

28:        Compute per-sample loss  $\ell_i \leftarrow \|\epsilon_i - \hat{\epsilon}_\theta\|^2$ 
29:         $L \leftarrow L + w_i \cdot \ell_i$ 
30:      end for
31:      Update  $\theta$  by one AdamW step with  $\eta$  and  $(\beta_1, \beta_2)$  using  $\nabla_\theta L$  clipped by GlobalNorm with  $C_{\text{clip}}$ 
32:    end for
33:  until convergence

```

---

### 3.2 Sparse Reward-Weighted Regression

Sparse Reward-Weighted Regression (RWR<sub>sparse</sub>) is similar to RWR. We consider a simplified weighting scheme that uses binary weights,

$$w_{\text{sparse}}(x_0, c) = 1_{\{r(x_0, c) \geq C\}},$$

where  $C$  is a reward threshold determining which samples are used for training.

We have the following pseudo-code for RWR<sub>sparse</sub> (a policy gradient method):

---

**Algorithm 2** Denoising Diffusion with Sparse Reward-Weighted Regression (RWR<sub>sparse</sub>)

---

```

1: Input: initial diffusion parameters  $\theta$ ; reward function  $r(x, c)$ ; prompt distribution  $p(c)$ ; timestep distribution  $p(t)$ .
2: Hyperparameters: reward temperature  $\beta$ ; Sparse Cap  $C$ ; number of gradient updates per iteration  $N_{\text{updates}}$ ; number
   of diffusion steps  $T$ ; CFG weight  $W_c$ ; number of samples per iteration  $N_{\text{samples}}$ ; number of inference steps  $N_{\text{inf}}$ ;
   batch size  $B$ ; noise schedule endpoints  $\beta_{\text{start}}, \beta_{\text{end}}$ ; learning rate  $\eta$ ; Adam parameters  $(\beta_1, \beta_2)$ ; Adam numerical
   constant  $\varepsilon$ ; weight decay  $\lambda$ ; gradient clipping norm  $C_{\text{clip}}$ .
3: repeat
4:   Initialize empty dataset  $\mathcal{D}$ 
5:   Use a diffusion scheduler with  $(T, \beta_{\text{start}}, \beta_{\text{end}})$  to generate the cumulative factors  $\{\bar{\alpha}_t\}_{t=0}^{T-1}$ 
6:   for  $i = 1$  to  $N_{\text{samples}}$  do
7:     Sample prompt  $c_i \sim p(c)$ 
8:     Run policy  $\pi_\theta(\cdot | c_i, w)$  to get a trajectory with  $N_{\text{inf}}$  time steps, the reward in each step is  $r_{\text{raw}}^{(i)} \leftarrow r(x_0^{(i)}, c_i)$ 
9:     Collect the terminal  $(c_i, x_0^{(i)}, r_{\text{raw}}^{(i)})$  to  $\mathcal{D}$ 
10:    end for
11:    for all  $(c_i, x_0^{(i)}, r_{\text{raw}}^{(i)}) \in \mathcal{D}$  do
12:       $w_i \leftarrow \mathbf{1}\{r_{\text{raw}}^{(i)} \geq C\}$ 
13:    end for
14:    for  $u = 1$  to  $N_{\text{updates}}$  do
15:      Sample a minibatch  $\mathcal{B} \subset \mathcal{D}$  of size  $B$ 
16:       $L \leftarrow 0$ 
17:      for all  $i \in \mathcal{B}$  do
18:        Sample timestep  $t_i \sim p(t)$  over  $\{0, \dots, T - 1\}$ 
19:        Sample noise  $\epsilon_i \sim \mathcal{N}(0, I)$ 
20:        Construct noised latent  $x_{t_i}^{(i)} \leftarrow \sqrt{\bar{\alpha}_{t_i}} x_0^{(i)} + \sqrt{1 - \bar{\alpha}_{t_i}} \epsilon_i$ 
21:        Compute conditional prediction  $\hat{\epsilon}_\theta^{\text{cond}} \leftarrow \epsilon_\theta(x_{t_i}^{(i)}, c_i, t_i)$ 
22:        Compute unconditional prediction  $\hat{\epsilon}_\theta^{\text{uncond}} \leftarrow \epsilon_\theta(x_{t_i}^{(i)}, \emptyset, t_i)$ 
23:        Apply classifier-free guidance:

$$\hat{\epsilon}_\theta \leftarrow \hat{\epsilon}_\theta^{\text{uncond}} + W_c (\hat{\epsilon}_\theta^{\text{cond}} - \hat{\epsilon}_\theta^{\text{uncond}})$$

24:        Compute per-sample loss  $\ell_i \leftarrow \|\epsilon_i - \hat{\epsilon}_\theta\|^2$ 
25:         $L \leftarrow L + w_i \cdot \ell_i$ 
26:      end for
27:      Update  $\theta$  by one AdamW step with  $\eta$  and  $(\beta_1, \beta_2)$  using  $\nabla_\theta L$  clipped by GlobalNorm with  $C_{\text{clip}}$ 
28:    end for
29:  until convergence

```

---

### 3.3 DDPO<sub>SF</sub>

DDPO is a form of Policy gradient estimation. With access to likelihoods and likelihood gradients, we can make direct Monte Carlo estimates of  $\nabla_\theta J_{\text{DDRL}}$ . DDPO alternates collecting denoising trajectories  $\{x_T, x_{T-1}, \dots, x_0\}$  via sampling and updating parameters via gradient descent.

The first variant of DDPO, which we call DDPO<sub>SF</sub>, uses the score function policy gradient estimator, also known as the likelihood ratio method or REINFORCE [3]:

$$\nabla_\theta J_{\text{DDRL}} = \mathbb{E} \left[ \sum_{t=0}^T \nabla_\theta \log p_\theta(x_{t-1} | x_t, c) r(x_0, c) \right] \quad (\text{DDPO}_{\text{SF}}) \quad (1)$$

where the expectation is taken over denoising trajectories generated by the current parameters  $\theta$ .

We have the following pesudo-code for DDPO<sub>SF</sub>:

---

**Algorithm 3** Denoising Diffusion Policy Optimization with Score-Function Estimator (DDPO<sub>SF</sub>)

---

1: **Input:** initial diffusion parameters  $\theta$ ; reward function  $r(x_0, c)$ ; prompt distribution  $p(c)$ .  
 2: **Hyperparameters:** number of denoising steps  $T$ ; number of trajectories per iteration  $N_{\text{samples}}$ ; number of gradient updates per iteration  $N_{\text{updates}}$ ; classifier-free guidance weight  $W_c$ ; batch size  $B$ ; noise schedule endpoints  $\beta_{\text{start}}, \beta_{\text{end}}$ ; learning rate  $\eta$ ; AdamW parameters  $(\beta_1, \beta_2, \varepsilon, \lambda)$ ; gradient clipping norm  $C_{\text{clip}}$ .  
 3: **repeat**  
 4:   Initialize empty dataset  $\mathcal{D} \leftarrow \emptyset$   
 5:   Use a diffusion scheduler with  $(T, \beta_{\text{start}}, \beta_{\text{end}})$  to generate variances  $\{\sigma_t^2\}_{t=1}^T$  and  $\{\beta_t\}_{t=1}^T$   
 6:   **for**  $i = 1$  to  $N_{\text{samples}}$  **do**  
 7:     Sample prompt  $c_i \sim p(c)$   
 8:     Sample initial latent  $x_T^{(i)} \sim \mathcal{N}(0, I)$   
 9:     Initialize log-likelihood accumulator  $\ell^{(i)} \leftarrow 0$   
 10:    **for**  $t = T, T-1, \dots, 1$  **do**  
 11:      Compute conditional prediction  $\hat{\epsilon}_\theta^{\text{cond}} \leftarrow \epsilon_\theta(x_t^{(i)}, c_i, t)$   
 12:      Compute unconditional prediction  $\hat{\epsilon}_\theta^{\text{uncond}} \leftarrow \epsilon_\theta(x_t^{(i)}, \emptyset, t)$   
 13:      Apply classifier-free guidance:  

$$\hat{\epsilon}_\theta \leftarrow \hat{\epsilon}_\theta^{\text{uncond}} + W_c(\hat{\epsilon}_\theta^{\text{cond}} - \hat{\epsilon}_\theta^{\text{uncond}})$$
  
 14:      Compute DDPM mean using the scheduler:  

$$\mu_t^{(i)} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( x_t^{(i)} - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \hat{\epsilon}_\theta \right)$$
  
 15:      Sample next latent  $x_{t-1}^{(i)} \sim \mathcal{N}(\mu_t^{(i)}, \sigma_t^2 I)$   
 16:      Accumulate log-likelihood term  

$$\ell^{(i)} \leftarrow \ell^{(i)} + \log \mathcal{N}(x_{t-1}^{(i)}; \mu_t^{(i)}, \sigma_t^2 I)$$
  
 17:    **end for**  
 18:    Let  $x_0^{(i)}$  be the final latent of the trajectory  
 19:    Compute terminal reward  $r_{\text{raw}}^{(i)} \leftarrow r(x_0^{(i)}, c_i)$   
 20:    Store  $(\ell^{(i)}, r_{\text{raw}}^{(i)}, c_i)$  in  $\mathcal{D}$   
 21:   **end for**  
 22:   **for**  $u = 1$  to  $N_{\text{updates}}$  **do**  
 23:     Sample a minibatch  $\mathcal{B} \subset \mathcal{D}$  of size  $B$   
 24:      $L \leftarrow 0$   
 25:     **for all**  $i \in \mathcal{B}$  **do**  
 26:        $L \leftarrow L - \tilde{r}_i \cdot \ell^{(i)}$   
 27:     **end for**  
 28:     Update  $\theta$  by one AdamW step with parameters  $(\eta, \beta_1, \beta_2, \varepsilon, \lambda)$ , using  $\nabla_\theta L$  clipped by global norm  $C_{\text{clip}}$   
 29:   **end for**  
 30: **until** convergence

---

### 3.4 DDPO<sub>IS</sub>

To perform multiple steps of optimization, we may use an importance sampling estimator:

$$\nabla_\theta J_{\text{DDRL}} = \mathbb{E} \left[ \sum_{t=0}^T \frac{p_\theta(x_{t-1} | x_t, c)}{p_{\theta_{\text{old}}}(x_{t-1} | x_t, c)} \nabla_\theta \log p_\theta(x_{t-1} | x_t, c) r(x_0, c) \right] \quad (\text{DDPO}_{\text{IS}}) \quad (2)$$

where the expectation is taken over denoising trajectories generated by the parameters  $\theta_{\text{old}}$ . This estimator becomes inaccurate if  $p_\theta$  deviates too far from  $p_{\theta_{\text{old}}}$ , which can be addressed using trust regions to constrain the size of the update. In practice, we implement the trust region via clipping.

---

**Algorithm 4** Denoising Diffusion Policy Optimization with Importance Sampling (DDPO<sub>IS</sub>)

---

1: **Input:** initial diffusion parameters  $\theta$ ; reward function  $r(x_0, c)$ ; prompt distribution  $p(c)$ .  
 2: **State:** behaviour parameters  $\theta_{\text{beh}} \leftarrow \theta$ .  
 3: **Hyperparameters:** number of denoising steps  $T$ ; number of trajectories per iteration  $N_{\text{samples}}$ ; number of gradient updates per iteration  $N_{\text{updates}}$ ; classifier-free guidance weight  $W_c$ ; batch size  $B$ ; noise schedule endpoints  $\beta_{\text{start}}, \beta_{\text{end}}$ ; learning rate  $\eta$ ; AdamW parameters  $(\beta_1, \beta_2, \varepsilon, \lambda)$ ; gradient clipping norm  $C_{\text{clip}}$ .  
 4: **repeat**  
 5:   Initialize empty dataset  $\mathcal{D} \leftarrow \emptyset$   
 6:   Set behaviour parameters to the current policy snapshot:  $\theta_{\text{beh}} \leftarrow \theta$   
 7:   Use a diffusion scheduler with  $(T, \beta_{\text{start}}, \beta_{\text{end}})$  to generate variances  $\{\sigma_t^2\}_{t=1}^T$  and  $\{\beta_t\}_{t=1}^T$   
 8:   **for**  $i = 1$  to  $N_{\text{samples}}$  **do**  
 9:     Sample prompt  $c_i \sim p(c)$   
 10:    Sample initial latent  $x_T^{(i)} \sim \mathcal{N}(0, I)$   
 11:    Initialize behaviour log-likelihood accumulator  $\ell_{\text{beh}}^{(i)} \leftarrow 0$   
 12:    **for**  $t = T, T - 1, \dots, 1$  **do**  
 13:      Compute conditional prediction  $\hat{\epsilon}_{\theta_{\text{beh}}}^{\text{cond}} \leftarrow \epsilon_{\theta_{\text{beh}}}(x_t^{(i)}, c_i, t)$   
 14:      Compute unconditional prediction  $\hat{\epsilon}_{\theta_{\text{beh}}}^{\text{uncond}} \leftarrow \epsilon_{\theta_{\text{beh}}}(x_t^{(i)}, \emptyset, t)$   
 15:      Apply classifier-free guidance:  

$$\hat{\epsilon}_{\theta_{\text{beh}}} \leftarrow \hat{\epsilon}_{\theta_{\text{beh}}}^{\text{uncond}} + W_c (\hat{\epsilon}_{\theta_{\text{beh}}}^{\text{cond}} - \hat{\epsilon}_{\theta_{\text{beh}}}^{\text{uncond}})$$
  
 16:      Compute DDPM mean using the scheduler:  

$$\mu_t^{(i)} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( x_t^{(i)} - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon}_{\theta_{\text{beh}}} \right)$$
  
 17:      Sample next latent  $x_{t-1}^{(i)} \sim \mathcal{N}(\mu_t^{(i)}, \sigma_t^2 I)$   
 18:      Accumulate behaviour log-likelihood term  

$$\ell_{\text{beh}}^{(i)} \leftarrow \ell_{\text{beh}}^{(i)} + \log \mathcal{N}(x_{t-1}^{(i)}; \mu_t^{(i)}, \sigma_t^2 I)$$
  
 19:    **end for**  
 20:    Let  $x_0^{(i)}$  be the final latent of the trajectory  
 21:    Compute terminal reward  $r_{\text{raw}}^{(i)} \leftarrow r(x_0^{(i)}, c_i)$   
 22:    Store trajectory  $\tau_i = \{x_t^{(i)}\}_{t=0}^T$  and stats  $(\tau_i, c_i, \ell_{\text{beh}}^{(i)}, r_{\text{raw}}^{(i)})$  in  $\mathcal{D}$   
 23: **end for**  
 24: **for**  $u = 1$  to  $N_{\text{updates}}$  **do**  
 25:    Sample a minibatch  $\mathcal{B} \subset \mathcal{D}$  of size  $B$   
 26:     $L \leftarrow 0$   
 27:    Compute any normalized rewards  $\{\tilde{r}_i\}_{i \in \mathcal{B}}$  from  $\{r_{\text{raw}}^{(i)}\}$   
 28:    **for all**  $i \in \mathcal{B}$  **do**  
 29:      Recompute trajectory log-likelihood under current parameters:  $\ell_{\theta}^{(i)} \leftarrow \log p_{\theta}(\tau_i)$   
 30:      Compute importance weight  

$$\tilde{w}_i \leftarrow \exp(\ell_{\theta}^{(i)} - \ell_{\text{beh}}^{(i)})$$
  
 31:      Accumulate loss with IS-corrected score-function term:  

$$L \leftarrow L - (\tilde{w}_i \tilde{r}_i) \cdot \ell_{\theta}^{(i)}$$
  
 32:    **end for**  
 33:    Update  $\theta$  by one AdamW step with parameters  $(\eta, \beta_1, \beta_2, \varepsilon, \lambda)$ , using  $\nabla_{\theta} L$  clipped by global norm  $C_{\text{clip}}$   
 34: **end for**  
 35: **until** convergence

---

### 3.5 DDPO<sub>ISKL</sub>

Here, we propose our own improved algorithm: DDPO<sub>ISKL</sub>. Our DDPO<sub>ISKL</sub> is based on DDPO<sub>IS</sub>, but we add a KL-Divergence regularization term between a pretrained diffusion model  $\theta_{\text{ref}}$  as reference, and our model in training  $\theta$ .

---

**Algorithm 5** Denoising Diffusion Policy Optimization with Importance Sampling and KL (DDPOISKL)

---

1: **Input:** initial diffusion parameters  $\theta$ ; fixed pretrained reference parameters  $\theta_{\text{ref}}$ ; reward function  $r(x_0, c)$ ; prompt distribution  $p(c)$ .  
2: **Hyperparameters:** number of denoising steps  $T$ ; number of trajectories per iteration  $N_{\text{samples}}$ ; number of gradient updates per iteration  $N_{\text{updates}}$ ; classifier-free guidance weight  $W_c$ ; batch size  $B$ ; noise schedule endpoints  $\beta_{\text{start}}, \beta_{\text{end}}$ ; learning rate  $\eta$ ; AdamW parameters  $(\beta_1, \beta_2, \varepsilon, \lambda)$ ; gradient clipping norm  $C_{\text{clip}}$ ; KL coefficient  $\lambda_{\text{KL}}$ .

3: **repeat**

4:   Set behaviour parameters to the current policy snapshot:  $\theta_{\text{beh}} \leftarrow \theta$   
5:   Initialize empty dataset  $\mathcal{D} \leftarrow \emptyset$   
6:   Use a diffusion scheduler with  $(T, \beta_{\text{start}}, \beta_{\text{end}})$  to generate variances  $\{\sigma_t^2\}_{t=1}^T$  and  $\{\beta_t\}_{t=1}^T$   
7:   **for**  $i = 1$  to  $N_{\text{samples}}$  **do**  
8:     Sample prompt  $c_i \sim p(c)$   
9:     Sample initial latent  $x_T^{(i)} \sim \mathcal{N}(0, I)$   
10:    Initialize behaviour log-likelihood accumulator  $\ell_{\text{beh}}^{(i)} \leftarrow 0$   
11:    **for**  $t = T, T-1, \dots, 1$  **do**  
12:      Compute conditional prediction  $\hat{\epsilon}_{\theta_{\text{beh}}}^{\text{cond}} \leftarrow \epsilon_{\theta_{\text{beh}}}(x_t^{(i)}, c_i, t)$   
13:      Compute unconditional prediction  $\hat{\epsilon}_{\theta_{\text{beh}}}^{\text{uncond}} \leftarrow \epsilon_{\theta_{\text{beh}}}(x_t^{(i)}, \emptyset, t)$   
14:      Apply classifier-free guidance:  

$$\hat{\epsilon}_{\theta_{\text{beh}}} \leftarrow \hat{\epsilon}_{\theta_{\text{beh}}}^{\text{uncond}} + W_c (\hat{\epsilon}_{\theta_{\text{beh}}}^{\text{cond}} - \hat{\epsilon}_{\theta_{\text{beh}}}^{\text{uncond}})$$
  
15:      Compute DDPM mean using the scheduler:  

$$\mu_t^{(i)} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( x_t^{(i)} - \frac{\beta_t}{\sqrt{1-\alpha_t}} \hat{\epsilon}_{\theta_{\text{beh}}} \right)$$
  
16:      Sample next latent  $x_{t-1}^{(i)} \sim \mathcal{N}(\mu_t^{(i)}, \sigma_t^2 I)$   
17:      Accumulate behaviour log-likelihood term  

$$\ell_{\text{beh}}^{(i)} \leftarrow \ell_{\text{beh}}^{(i)} + \log \mathcal{N}(x_{t-1}^{(i)}; \mu_t^{(i)}, \sigma_t^2 I)$$
  
18:    **end for**  
19:    Let  $x_0^{(i)}$  be the final latent of the trajectory  
20:    Compute terminal reward  $r_{\text{raw}}^{(i)} \leftarrow r(x_0^{(i)}, c_i)$   
21:    Store trajectory  $\tau_i = \{x_t^{(i)}\}_{t=0}^T$  and stats  $(\tau_i, c_i, \ell_{\text{beh}}^{(i)}, r_{\text{raw}}^{(i)})$  in  $\mathcal{D}$   
22: **end for**  
23: **for**  $u = 1$  to  $N_{\text{updates}}$  **do**  
24:    Sample a minibatch  $\mathcal{B} \subset \mathcal{D}$  of size  $B$   
25:     $L \leftarrow 0$   
26:    Compute any normalized rewards  $\{\tilde{r}_i\}_{i \in \mathcal{B}}$  from  $\{r_{\text{raw}}^{(i)}\}$   
27:    **for all**  $i \in \mathcal{B}$  **do**  
28:      Recompute trajectory log-likelihood under current parameters:  $\ell_{\theta}^{(i)} \leftarrow \log p_{\theta}(\tau_i)$   
29:      Compute reference log-likelihood under fixed pretrained parameters:  $\ell_{\text{ref}}^{(i)} \leftarrow \log p_{\theta_{\text{ref}}}(\tau_i)$   
30:      Compute importance weight  

$$\tilde{w}_i \leftarrow \exp(\ell_{\theta}^{(i)} - \ell_{\text{beh}}^{(i)})$$
  
31:      Accumulate loss with IS-corrected score-function term and KL penalty:  

$$L \leftarrow L - (\tilde{w}_i \tilde{r}_i) \cdot \ell_{\theta}^{(i)} + \frac{\lambda_{\text{KL}}}{2} (\ell_{\theta}^{(i)} - \ell_{\text{ref}}^{(i)})^2$$
  
32: **end for**  
33:    Update  $\theta$  by one AdamW step with parameters  $(\eta, \beta_1, \beta_2, \varepsilon, \lambda)$ , using  $\nabla_{\theta} L$  clipped by global norm  $C_{\text{clip}}$   
34: **end for**  
35: **until** convergence

---

## 4 Experiment and Results

### 4.1 Experimental Setup

We run the three algorithms :DDPO<sub>SF</sub>, DDPO<sub>IS</sub>, DDPO<sub>ISKL</sub> (RWR, RWR<sub>Sparse</sub> are omitted because they are not really a part of this paper, but from previous literature) described above across three different downstream tasks (Compressibility and Incompressibility, Aesthetic Quality).

For DDPO<sub>ISKL</sub>, we also choose two sets of hyperparameter for configuration of model training, one with fairly-optimal results and one with sub-optimal results. We also plot the reward and trajectories (against the number of iterations) for each configuration.

The following table 1 is a full documentation of our configuration of hyperparameters for the model to work well:

Table 1: A set of Hyperparameter that work well

	<b>DDPO<sub>IS</sub></b>	<b>DDPO<sub>SF</sub></b>	<b>DDPO<sub>ISKL</sub></b>
<b>Diffusion</b>			
Denoising steps ( $T$ )	50	50	50
Guidance weight ( $w$ )	5.0	5.0	5.0
<b>Optimization</b>			
Optimizer	AdamW	AdamW	AdamW
Learning rate	$3 \times 10^{-4}$	$3 \times 10^{-4}$	$3 \times 10^{-4}$
Weight decay	$1 \times 10^{-4}$	$1 \times 10^{-4}$	$1 \times 10^{-4}$
$\beta_1$	0.9	0.9	0.9
$\beta_2$	0.999	0.999	0.999
$\epsilon$	$1 \times 10^{-8}$	$1 \times 10^{-8}$	$1 \times 10^{-8}$
Gradient clip norm	1.0	1.0	1.0
<b>DDPO</b>			
Batch size	16	16	16
Samples per iteration	8	16	8
Gradient updates per iteration	2	4	2
Clip range	$1 \times 10^{-4}$	-	$1 \times 10^{-4}$

For sub-optimal hyperparamters (Meaning bad hyperparameters), the only change we make compared to the hyperparameters above is we change the learning rate to  $3 \times 10^{-3}$  (from  $3 \times 10^{-4}$  ).

### 4.2 Experiment Results

We report the results of the experiments of the three algorithms DDPO<sub>SF</sub>, DDPO<sub>IS</sub>, and DDPO<sub>ISKL</sub> across some selected tasks in the following format:

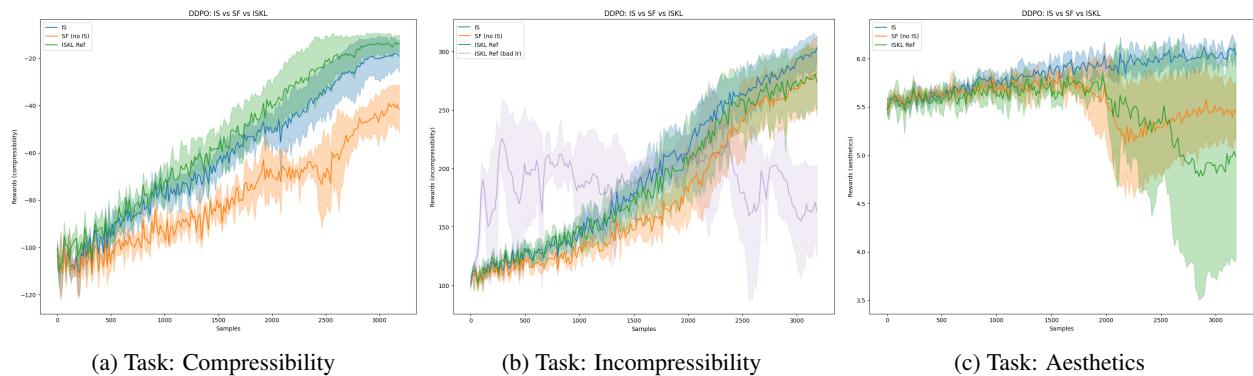


Figure 3: Rewards of three algorithms v.s. number of samples for downstream tasks, reported with mean value and standard deviation for 5 runs per algo

Figure 3 shows the plots of rewards against the number of samples. Each plot shows the rewards against samples for each one of three tasks. Each plot contains three curves that are the rewards plotted against number of samples for the three algorithms. The rewards are averaged over 5 runs with shading area representing the standard deviation from the mean. This is with good hyperparameters from above.

Figure 3b shows the plots of rewards against the iterations of algorithms, but with bad hyperparameters (learning rate  $3 \times 10^{-3}$ ) included. The algorithm chosen to demonstrate bad hyperparameters in DDPO<sub>ISKL</sub>.

For DDPO<sub>ISKL</sub>, we also report two groups of plots, each group is configured with a different set of hyperparameters, one working well, the other less optimal.

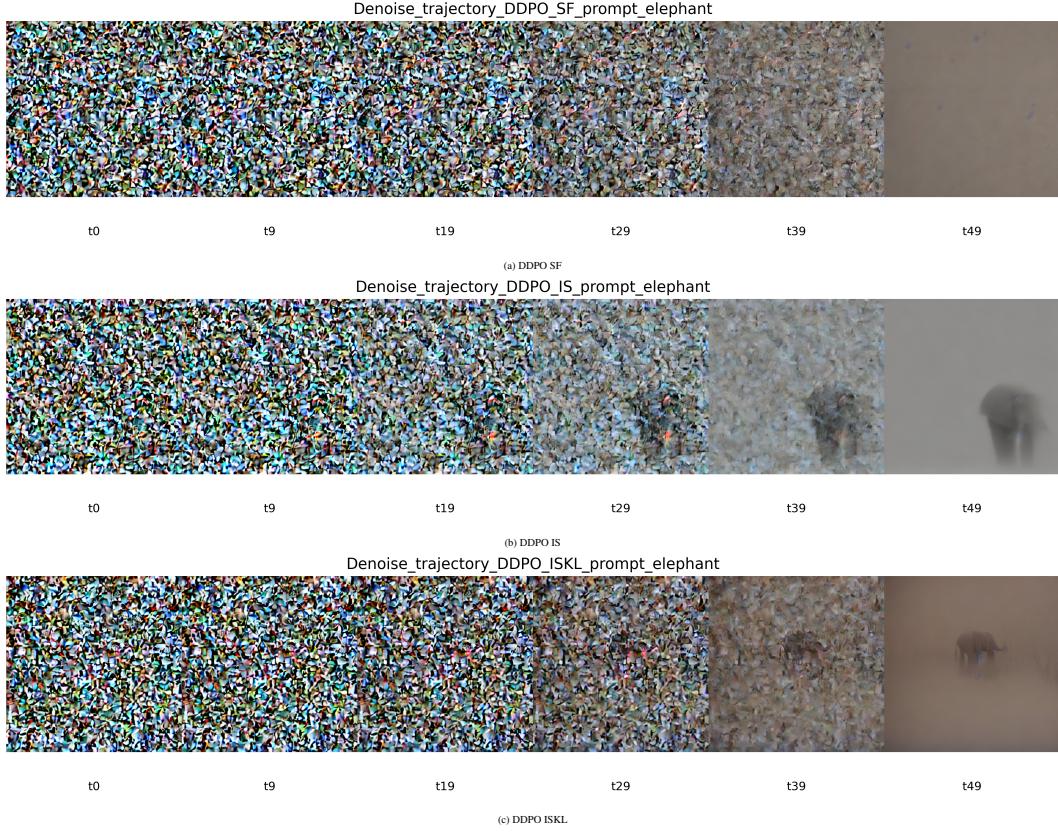


Figure 4: Denoising trajectories for three algorithms with prompt "elephant" with good parameters

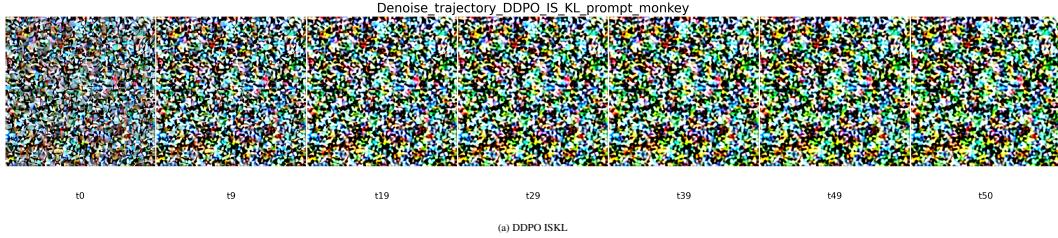


Figure 5: Some Denoising trajectories for DDPO<sub>ISKL</sub> with bad hyperparameters

Figure 4 shows the plots of sampled trajectories with the learned policy against time steps for the selected task Compressibility. All three algorithms are shown. This is with good hyperparameters from above.

Firgure 5 shows the plots of sampled trajectories with the learned policy against time steps for the selected task Compressibility, but with bad hyperparameters (change the learning rate to  $3 \times 10^{-3}$ ). The algorithm chosen to demonstrate the hyperparameters is DDPO<sub>ISKL</sub>.

## 5 Conclusion

We have successfully replicated the paper: *Training Diffusion Models with Reinforcement Learning*; Training diffusion models different reinforcement learning algorithms across various downstream tasks. We also proposed our own improvement to the algorithm DDPO<sub>IS</sub>, hence formulating our own algorithm DDPO<sub>ISKL</sub>, achieving faster convergence and avoiding overoptimization in some tasks.

For all three tasks we observed that with Important Sampling, DDPO converges faster and show more stable training and less variance. The effect can be attributed to 2 reasons: 1) Off-policy allows more frequent gradient update from one round of data sampling 2) important sampling with clipping make the gradient less noisy thus resulting in more stable training process.

The added KL constraint with respect to a reference Stable Diffusion model is intended to keep the learned distribution close to the pretrained one, preventing the policy from collapsing into nonsense generations such as the SF algorithm's output shown here. It has a high score of compressibility, but completely loses the text-to-image alignment. Empirically, combining KL with importance sampling (IS) leads to noticeably better visual quality. Intuitively, the KL penalty regularizes the updated distribution and discourages overly aggressive parameter updates that would otherwise produce meaningless noise. As a result, the model is guided toward more desirable sampling trajectories, reducing the failure of completely losing alignment for pursuing high rewards. On tasks that are harder to optimize, such as aesthetics score in our case, the KL introduced extra complexity to gradient update. This could disturb the training when gradient is unstable due to our small batch size, where only DDPO IS has not suffered from collapsing during training.

## References

- [1] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning, 2024.
- [2] Zejian Li, Chenye Meng, Yize Li, Ling Yang, Shengyuan Zhang, Jiarui Ma, Jiayi Li, Guang Yang, Changyuan Yang, Zhiyuan Yang, Jinxiang Chang, and Lingyun Sun. Laion-sg: An enhanced large-scale dataset for training complex image-text models with structural annotations, 2024.
- [3] Junzi Zhang, Jongho Kim, Brendan O'Donoghue, and Stephen Boyd. Sample efficient reinforcement learning with reinforce, 2020.