

# Diffusion Policy Policy Optimization

Allen Z. Ren<sup>1</sup>, Justin Lidard<sup>1</sup>, Lars L. Ankile<sup>2</sup>, Anthony Simeonov<sup>2</sup>,  
 Pulkit Agrawal<sup>2</sup>, Anirudha Majumdar<sup>1</sup>, Benjamin Burchfiel<sup>3</sup>, Hongkai Dai<sup>3</sup>, Max Simchowitz<sup>2,4</sup>

<sup>1</sup>Princeton University <sup>2</sup>Massachusetts Institute of Technology  
<sup>3</sup>Toyota Research Institute <sup>4</sup>Carnegie Mellon University

## Abstract

We introduce *Diffusion Policy Policy Optimization*, **DPPPO**, an algorithmic framework including best practices for fine-tuning diffusion-based policies (e.g. Diffusion Policy [15]) in continuous control and robot learning tasks using the policy gradient (PG) method from reinforcement learning (RL). PG methods are ubiquitous in training RL policies with other policy parameterizations; nevertheless, they had been conjectured to be less efficient for diffusion-based policies. Surprisingly, we show that **DPPPO** achieves the strongest overall performance and efficiency for fine-tuning in common benchmarks compared to other RL methods for diffusion-based policies and also compared to PG fine-tuning of other policy parameterizations. Through experimental investigation, we find that **DPPPO** takes advantage of unique synergies between RL fine-tuning and the diffusion parameterization, leading to structured and on-manifold exploration, stable training, and strong policy robustness. We further demonstrate the strengths of **DPPPO** in a range of realistic settings, including simulated robotic tasks with pixel observations, and via zero-shot deployment of simulation-trained policies on robot hardware in a long-horizon, multi-stage manipulation task. Website with code: [diffusion-ppo.github.io](https://diffusion-ppo.github.io).

## 1 Introduction

Large-scale pre-training with additional fine-tuning has become a ubiquitous pipeline in the development of language and image foundation models [9, 60, 51, 67]. Though behavior cloning with expert data [56] is rapidly emerging as dominant paradigm for pre-training *robot policies* [22, 23, 88, 40, 25], their performance can be suboptimal [50] due to expert data being suboptimal or expert data exhibiting limited coverage of possible environment conditions. As robot policies entail interaction with their environment, reinforcement learning (RL) [75] is a natural candidate for further optimizing their performance beyond the limits of demonstration data. However, RL fine-tuning can be nuanced for pre-trained policies parameterized as diffusion models [29], which have emerged as a leading parameterization for action policies [15, 63, 52], due in large part to their high training stability and ability to represent complex distributions [65, 57, 39, 30].

**Contribution 1 (DPPPO).** We introduce *Diffusion Policy Policy Optimization* (**DPPPO**), a generic framework as well as a set of carefully chosen design decisions for fine-tuning a diffusion-based robot learning policy via popular policy gradient methods [76, 70] in reinforcement learning.

The literature has already studied improving/fine-tuning diffusion-based policies (*Diffusion Policy*) using RL [59, 84, 27], and has applied policy gradient (PG) to fine-tuning non-interactive applications of diffusion models such as text-to-image generation [6, 16, 20]. Yet PG methods have been believed to be inefficient in training Diffusion Policy for continuous control tasks [59, 86]. On the contrary, we show that for a Diffusion Policy pre-trained from expert demonstrations, our methodology for *fine-tuning* via PG updates yields robust, high-performing policies with favorable training behavior.

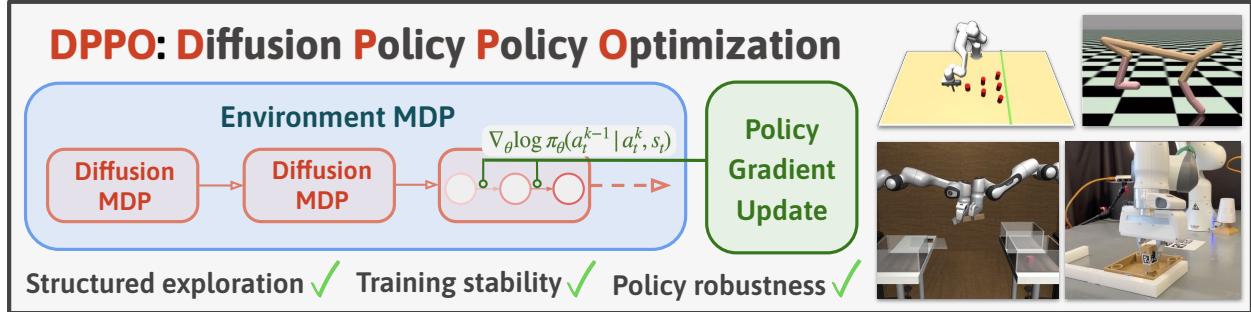


Figure 1: We introduce **DPPO**, *Diffusion Policy Policy Optimization*, that fine-tunes pre-trained Diffusion Policy using policy gradient. **DPPO** affords structured exploration and training stability during policy fine-tuning, and the fine-tuned policy exhibits strong robustness and generalization. **DPPO** improves policy performance across benchmarks, including ones with pixel observations and with long-horizon rollouts that have been very challenging to solve using previous RL methods.

**Contribution 2 (Demonstration of DPPO’s Performance).** We show that for *fine-tuning* a pre-trained Diffusion Policy, **DPPO** yields consistent and marked improvements in training stability and final policy performance in comparison to a range of alternatives, including those based on off-policy Q-learning [84, 27, 86, 59] and weighted regression [53, 55, 37], as well as common policy parameterizations such as Gaussian and Gaussian Mixture models.

The above finding might be surprising because PG methods do not appear to take advantage of the unique capabilities of diffusion sampling (e.g., guidance [35, 2]). Through careful investigative experimentation, however, we find a **unique synergy** between RL fine-tuning and diffusion-based policies.

**Contribution 3 (Understanding the mechanism of DPPO’s success).** We complement our results with numerous investigative experiments that provide insight into the mechanisms behind **DPPO**’s strong performance. Compared to other common policy parameterizations, we provide evidence that **DPPO** engages in *structured exploration* that takes better advantage of the “manifold” of training data, and finds policies that exhibit greater robustness to perturbation.

Through ablations, we further show that our design decisions overcome the speculated limitation of PG methods for fine-tuning Diffusion Policy. Finally, to justify the broad utility of **DPPO**, we verify its efficacy across both simulated and real environments, and in situations when either ground-truth states or pixels are given to the policy as input.

**Contribution 4 (Tackling challenging robotic tasks and settings).** We demonstrate the effectiveness of **DPPO** in different challenging robotic and control settings, including ones with pixel observations and long-horizon manipulation tasks with sparse reward. We also deploy a policy trained in simulation via **DPPO** on real hardware in zero-shot, which exhibits a remarkably small sim-to-real gap compared to the baseline.

**Potential impact beyond robotics.** **DPPO** is a generic framework that can be potentially applied to fine-tuning diffusion-based models in sequential interactive settings beyond robotics. These include: extending diffusion-based text-to-image generation [6, 16] to a multi-turn interactive setting with human feedback; drug design/discovery applications [46, 33] with policy search on the molecular level in feedback with simulators (in the spirit of prior non-diffusion-based drug discovery with RL [58]); and the adaptation of diffusion-based language modeling [68, 45] to interactive (e.g. with human feedback [51]), problem-solving and planning tasks.

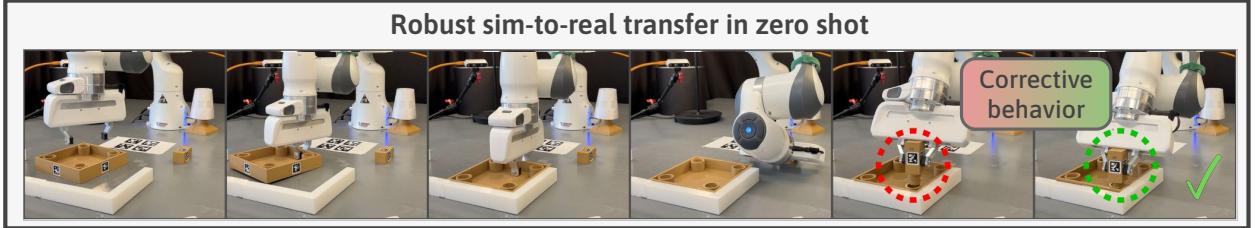


Figure 2: **DPO** solves challenging long-horizon manipulation tasks from FURNITURE-BENCH ([28]), and allows robust sim-to-real transfer without using any real data (see [Section 5.3](#)).

## 1.1 Overview of Approach

Like prior work [6, 16], our approach unrolls the denoising diffusion process into an MDP in which action likelihoods are explicit; from this, we construct a *two-layer MDP* whose outer layer corresponds to the environment MDP and inner layer to the denoising MDP. Given the challenges of RL training, we also present a number of design decisions tailored to DPPO that are crucial for enabling **DPPO**’s performance in challenging robotic settings, discussed in **Contributions 2** and **4**, including:

- We apply *Proximal Policy Optimization* (PPO) [70] to the two-layer MDP. We show how to efficiently estimate the advantage function for the PPO update.
- We show that for fine-tuning tasks, it often suffices to fine-tune only the last few steps of the denoising process, or fine-tune the Denoising Diffusion Implicit Model (DDIM) [72] instead.
- We propose modifications to the diffusion noise schedule to ensure stable training and adequate exploration, whilst also leveraging the natural stochasticity of diffusion models.

Our method is formally detailed in [Section 4](#). As noted above, we conduct extensive experiments demonstrating both the success of **DPPO** in [Section 5](#) and potential explanations thereof in [Section 6](#).

## 2 Related Work

**Policy optimization and its application to robotics.** Policy optimization methods update an explicit representation of an RL policy — typically parameterized by a neural network — by taking gradients through action likelihoods. Following the seminal policy gradient (PG) method [85, 76], there have been a range of algorithms that further improve training stability and sample efficiency such as DDPG [43] and PPO [69]. PG methods have been broadly effective in training robot policies [3, 34, 38, 13], largely due to their training stability with high-dimensional continuous action spaces, as well as their favorable scaling with parallelized simulated environments. Given the challenges of from-scratch exploration in long horizon tasks, PG has seen great success in fine-tuning a baseline policy trained from demonstrations [61, 78, 54].

**Learning and improving diffusion-based policies.** Diffusion-based policies [15, 63, 4, 87, 83, 74, 52] have shown recent success in robotics and decision making applications. Most typically, these policies are trained from human demonstrations through a supervised objective, and enjoy both high training stability and strong performance in modeling complex and multi-modal trajectory distributions.

As demonstration data are often limited and/or suboptimal, there have been many approaches proposed to improve the performance of diffusion-based policies. One popular approach has been to guide the diffusion denoising process using objectives such as reward signal or goal conditioning [35, 2, 42, 80, 11]. More

recent work has explored the possibilities of techniques including Q-learning and weighted regression, either from purely offline estimation [12, 84, 17], and/or with online interaction [37, 27, 59, 86].

**Policy gradient through diffusion models.** RL techniques have been used to fine-tune diffusion models such as ones for text-to-image generation [19, 20, 6, 16, 81]. Clark et al. [16] treat the denoising process as an MDP, to which they apply policy gradient updates, and Black et al. [6] apply the PPO update to this formalism. We build upon these earlier findings by embedding the denoising MDP into the environmental MDP of the dynamics in control tasks, forming a two-layer ‘‘Diffusion Policy MDP’’. Though Psenka et al. [59] have already shown how PG can be taken through Diffusion Policy by propagating PG through both MDPs, they conjecture that it is likely to be ineffective due to large action variance caused by the increased effective horizon induced from the denoising steps. Our results contravene this supposition for diffusion-based policies in the fine-tuning setting.

### 3 Preliminaries

**Markov Decision Process.** We consider *Markov Decision Process* (MDP)<sup>1</sup>  $\mathcal{M}_{\text{ENV}} := (\mathcal{S}, \mathcal{A}, P_0, P, R)$  with states  $s \in \mathcal{S}$ , actions  $a \in \mathcal{A}$ , initial state distribution  $P_0$ , transition probabilities  $P$ , and reward  $R$ .

At each timestep  $t$ , the agent (e.g., robot) observes the state  $s_t \in \mathcal{S}$ , takes an action  $a_t \sim \pi(a_t | s_t) \in \mathcal{A}$ , transitions to the next state  $s_{t+1}$  according to  $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$  while receiving the reward  $R(s_t, a_t)$ <sup>2</sup>. Fixing the MDP  $\mathcal{M}_{\text{ENV}}$ , we let  $\mathbb{E}^\pi$  (resp.  $\mathbb{P}^\pi$ ) denote the expectation (resp. probability distribution) over trajectories,  $\tau = (s_0, a_0, \dots, s_T, a_T)$  with length  $T + 1$ , with initial state distribution  $s_0 \sim P_0$  and transition operator  $P$ . We aim to train a policy to optimize the cumulative reward, discounted by a function  $\gamma(\cdot)$ , such that the agent receives:

$$\mathcal{J}(\pi_\theta) = \mathbb{E}^{\pi_\theta, P_0} \left[ \sum_{t \geq 0} \gamma(t) R(s_t, a_t) \right]. \quad (3.1)$$

**Policy optimization.** The *policy gradient method* (e.g., REINFORCE [85]) allows for improving policy performance by approximating the gradient of this objective w.r.t. the policy parameters:

$$\nabla_\theta \mathcal{J}(\pi_\theta) = \mathbb{E}^{\pi_\theta, P_0} \left[ \sum_{t \geq 0} \nabla_\theta \log \pi_\theta(a_t | s_t) r_t(s_t, a_t) \right], \quad r_t(s_t, a_t) := \sum_{\tau \geq t} \gamma(\tau) R(s_\tau, a_\tau), \quad (3.2)$$

where  $r_t$  is the discounted cumulative future reward from time  $t$  (more generally,  $r_t$  can be replaced by a Q-function estimator [76]),  $\gamma$  is the discount factor that depends on the time-step, and  $\nabla_\theta \log \pi_\theta(a_t | s_t)$  denotes the gradient of the logarithm of the *likelihood* of  $a_t | s_t$ . To reduce the variance of the gradient estimation, a state-value function  $\hat{V}^{\pi_\theta}(s_t)$  can be learned to approximate  $\mathbb{E}[r_t]$ . The estimated advantage function  $\hat{A}^{\pi_\theta}(s_t, a_t) := r_t(s_t, a_t) - \hat{V}^{\pi_\theta}(s_t)$  substitutes  $r_t(s_t, a_t)$ .

**Diffusion models.** A denoising diffusion probabilistic model (DDPM) [49, 29, 71] represents a continuous-valued data distribution  $p(\cdot) = p(x^0)$  as the reverse denoising process of a forward noising process  $q(x^k | x^{k-1})$  that iteratively adds Gaussian noise to the data. The reverse process is parameterized by a neural network

---

<sup>1</sup>More generally, we can view our environment as a Partially Observed Markov Decision Process (POMDP) where the agent’s actions depend on observations  $o$  of the states  $s$  (e.g., action from pixels). Our implementation applies in this setting, but we omit additional observations from the formalism to avoid notional clutter.

<sup>2</sup>For simplicity, we overload  $R(\cdot, \cdot)$  to denote both the random variable reward and its distribution.

$\varepsilon_\theta(x_k, k)$ , predicting the added noise  $\varepsilon$  that converts  $x_0$  to  $x_k$  [29]. Sampling starts with a random sample  $x^K \sim \mathcal{N}(0, \mathbf{I})$  and iteratively generates the denoised sample:

$$x^{k-1} \sim p_\theta(x^{k-1} | x^k) := \mathcal{N}(x^{k-1}; \mu_k(x^k, \varepsilon_\theta(x^k, k)), \sigma_k^2 \mathbf{I}). \quad (3.3)$$

Above,  $\mu_k(\cdot)$  is a fixed function, independent of  $\theta$ , that maps  $x^k$  and predicted  $\varepsilon_\theta$  to the next mean, and  $\sigma_k^2$  is a variance term that abides by a fixed schedule from  $k = 1, \dots, K$ . We refer the reader to Chan [10] for an in-depth survey.

**Diffusion models as policies.** *Diffusion Policy* (DP; see [15]) is a policy  $\pi_\theta$  parameterized by a DDPM which takes in  $s$  as a conditioning argument, and parameterizes  $p_\theta(a^{k-1} | a^k, s)$  as in (3.3). DPs can be trained via behavior cloning by fitting the conditional noise prediction  $\varepsilon_\theta(a^k, s, k)$  to predict added noise. Notice that unlike more standard policy parameterizations such as unimodal Gaussian policies, DPs do not maintain an explicit likelihood of  $p_\theta(a^0 | s)$ . In this work, we adopt the common practice of training DPs to predict an **action chunk** — a sequence of actions a few time steps (denoted  $T_a$ ) into the future — to promote temporal consistency. For fair comparison, our non-diffusion baselines use the same chunk size.

## 4 DPPO: Diffusion Policy Policy Optimization

As noted in [Section 1](#) and [Section 2](#), there has been much attention devoted to fine-tuning Diffusion Policy for improved performance [86, 59], focusing primarily on off-policy Q-learning [84, 27] and/or weighted regression. **DPPO** takes a different approach: differentiate through action likelihood and apply a policy gradient (PG) update like (3.2).

### 4.1 A Two-Layer “Diffusion Policy MDP”

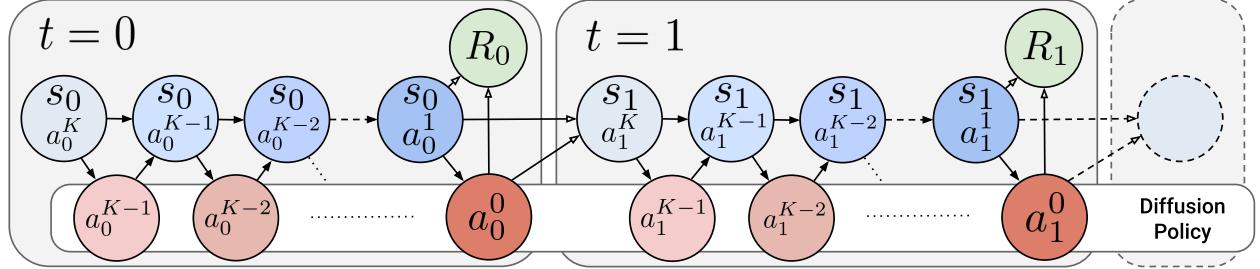


Figure 3: We treat the denoising process in Diffusion Policy as a MDP, and the whole environment episode can be considered as a chain of such MDPs. Now the entire chain (“Diffusion Policy MDP”,  $\mathcal{M}_{\text{DP}}$ ) involves a Gaussian likelihood at each (denoising) step and thus can be optimized with policy gradient. Blue circle denotes the state and red circle denotes the action in  $\mathcal{M}_{\text{DP}}$ .

As observed in [16, 6] and [59], a denoising process can be represented as a multi-step MDP in which policy likelihoods can be obtained directly. We extend this formalism by embedding the Diffusion MDP into the environmental MDP, obtaining a larger “Diffusion Policy MDP” denoted  $\mathcal{M}_{\text{DP}}$ , visualized in [Fig. 3](#). Below, we use the notation  $\delta$  to denote a Dirac distribution and  $\otimes$  to denote a product distribution.

Recall the environment MDP  $\mathcal{M}_{\text{ENV}} := (\mathcal{S}, \mathcal{A}, P_0, P, R)$  in [Section 3](#). The Diffusion MDP  $\mathcal{M}_{\text{DP}}$  uses indices  $\bar{t}(t, k) = tK + (K - k - 1)$  corresponding to  $(t, k)$ , which increases in  $t$  but (to keep the indexing conventions of diffusion) *decreases* lexicographically with  $K - 1 \geq k \geq 0$ . We write states, actions and

rewards as,

$$\bar{s}_{\bar{t}(t,k)} = (s_t, a_t^{k+1}), \quad \bar{a}_{\bar{t}(t,k)} = a_t^k, \quad \bar{R}_{\bar{t}(t,k)}(\bar{s}_{\bar{t}(t,k)}, \bar{a}_{\bar{t}(t,k)}) = \begin{cases} 0 & k > 0 \\ R(s_t, a_t^0) & k = 0 \end{cases},$$

where the bar-action at  $\bar{t}(t, k)$  is the action  $a_t^k$  after one denoising step. Reward is only given at times corresponding to when  $a_t^0$  is taken. The initial state distribution is  $\bar{P}^0 = P_0 \otimes \mathcal{N}(0, \mathbf{I})$ , corresponding to  $s_0 \sim P_0$  is the initial distribution from the environmental MDP and  $a_0^K \sim \mathcal{N}(0, \mathbf{I})$  independently. Finally, the transitions are

$$\bar{P}(\bar{s}_{\bar{t}+1} \mid \bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) = \begin{cases} (s_t, a_t^k) \sim \delta_{(s_t, a_t^k)} & \bar{t} = \bar{t}(t, k), k > 0 \\ (s_{t+1}, a_{t+1}^K) \sim P(s_{t+1} \mid s_t, a_t^0) \otimes \mathcal{N}(0, \mathbf{I}) & \bar{t} = \bar{t}(t, k), k = 0 \end{cases}.$$

That is, the transition moves the denoised action  $a_t^k$  at step  $\bar{t}(t, k)$  *into the next state* when  $k > 0$ , or otherwise progresses the environment MDP dynamics with  $k = 0$ . The pure noise  $a_t^K$  is considered part of the *environment* when transitioning at  $k = 0$ . In light of (3.3), the policy in  $\mathcal{M}_{\text{DP}}$  takes the form

$$\bar{\pi}_\theta(\bar{a}_{\bar{t}(t,k)} \mid \bar{s}_{\bar{t}(t,k)}) = \pi_\theta(a_t^k \mid a_t^{k+1}, s_t) = \mathcal{N}(a_t^k; \mu(a_t^{k+1}, \varepsilon_\theta(a_t^{k+1}, k+1, s_t)), \sigma_{k+1}^2 \mathbf{I}). \quad (4.1)$$

Fortunately, (4.1) is a *Gaussian likelihood*, which can be evaluated analytically and is amenable to the policy gradient updates (see also [59] for an alternative derivation):

$$\nabla_\theta \bar{\mathcal{J}}(\bar{\pi}_\theta) = \mathbb{E}^{\bar{\pi}_\theta, \bar{P}, \bar{P}^0} \left[ \sum_{\bar{t} \geq 0} \nabla_\theta \log \bar{\pi}_\theta(\bar{a}_{\bar{t}} \mid \bar{s}_{\bar{t}}) \bar{r}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) \right], \quad \bar{r}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) := \sum_{\tau \geq \bar{t}} \gamma(\tau) \bar{R}(\bar{s}_\tau, \bar{a}_\tau). \quad (4.2)$$

Evaluating the above involves sampling through the denoising process, which is the usual “forward pass” that samples actions in Diffusion Policy; as noted above, the initial state can be sampled from the environment via  $\bar{P}^0 = P_0 \otimes \mathcal{N}(0, \mathbf{I})$ , where  $P_0$  is from the environment MDP.

## 4.2 Instantiating DPPO with Proximal Policy Optimization

We apply Proximal Policy Optimization (PPO) [70, 18, 32, 1], a popular improvement of the vanilla policy gradient update.

**Definition 4.1** (Generalized PPO). Consider a general MDP. Given an advantage estimator  $\hat{A}(s, a)$ , the PPO update is given by [70] the sample approximation to

$$\nabla_\theta \mathbb{E}^{(s_t, a_t) \sim \pi_{\theta_{\text{old}}}} \min \left( \hat{A}^{\pi_{\theta_{\text{old}}}}(s_t, a_t) \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}, \hat{A}^{\pi_{\theta_{\text{old}}}}(s_t, a_t) \text{clip} \left( \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_{\text{old}}}(a_t \mid s_t)}, 1 - \varepsilon, 1 + \varepsilon \right) \right),$$

where  $\varepsilon$ , the clipping ratio, controls the maximum magnitude of policy change from the previous policy.

We instantiate PPO in our diffusion MDP with  $(s, a, t) \leftarrow (\bar{s}, \bar{a}, \bar{t})$ . Our advantage estimator takes a specific form that respects the two-level nature of the MDP: let  $\gamma_{\text{ENV}} \in (0, 1)$  be the environment discount and  $\gamma_{\text{DENOISE}} \in (0, 1)$  be the denoising discount. Consider the environment-discounted return:

$$\bar{r}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) := \sum_{t' \geq \bar{t}} \gamma_{\text{ENV}}^t \bar{r}(\bar{s}_{\bar{t}(t',0)}, \bar{a}_{\bar{t}(t',0)}), \quad \bar{t} = \bar{t}(t, k),$$

since  $\bar{R}(\bar{t}) = 0$  at  $k > 0$ . This fact also obviates the need of estimating the value at  $k > 1$  and allows us to use the following denoising-discounted advantage estimator<sup>3</sup>:

$$\hat{A}^{\pi_{\theta_{\text{old}}}}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) := \gamma_{\text{DENOISE}}^k \left( \bar{r}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) - \hat{V}^{\bar{\pi}_{\theta_{\text{old}}}}(\bar{s}_{\bar{t}(t,0)}) \right)$$

The denoising-discounting has the effect of downweighting the contribution of noisier steps (larger  $k$ ) to the policy gradient (see study in [Appendix B.1](#)). Lastly, we choose the value estimator to *only depend* on the “ $s$ ” component of  $\bar{s}$ :

$$\hat{V}^{\bar{\pi}_{\theta_{\text{old}}}}(\bar{s}_{\bar{t}(t,0)}) := \tilde{V}^{\bar{\pi}_{\theta_{\text{old}}}}(s_t),$$

which we find leads to more efficient and stable training compared to also estimating the value of applying the denoised action  $a_t^{k=1}$  (part of  $\bar{s}_{\bar{t}(t,0)}$ ) as shown in [Appendix B.1](#).

### 4.3 Best Practices for **DPO**

**Fine-tune only the last few denoising steps.** Diffusion Policy often uses up to  $K = 100$  denoising steps with DDPM to better capture the complex data distribution of expert demonstrations. With **DPO**, we can choose to fine-tune only a subset of the denoising steps instead, e.g., the last  $K'$  steps. Experimental results in [Appendix B.1](#) shows this speeds up **DPO** training and reduces GPU memory usage without sacrificing the final performance. Instead of fine-tuning the pre-trained model weights  $\theta$ , we make a copy  $\theta_{\text{FT}}$  —  $\theta$  is frozen and used for the early denoising steps, while  $\theta_{\text{FT}}$  is used for the last  $K'$  steps and updated with **DPO**.

**Fine-tune DDIM sampling.** Instead of fine-tuning all  $K$  or the last few steps of the DDPM, one can also apply Denoising Diffusion Implicit Model (DDIM) [72] during fine-tuning, which greatly reduces the number of sampling steps  $K^{\text{DDIM}} \ll K$ , e.g., as few as 5 steps, and thus potentially improves **DPO** efficiency as fewer steps are fine-tuned:

$$x^{k-1} \sim p_{\theta}^{\text{DDIM}}(x^{k-1}|x^k) := \mathcal{N}(x^{k-1}; \mu^{\text{DDIM}}(x^k, \varepsilon_{\theta}(x^k, k)), \eta \sigma_k^2 \mathbf{I}), \quad k = K^{\text{DDIM}}, \dots, 0. \quad (4.3)$$

Although DDIM is typically used as a deterministic sampler by setting  $\eta = 0$  in (4.3), we can use  $\eta > 0$  for fine-tuning in order to provide exploration noise and avoid calculating a Gaussian likelihood with a Dirac distribution. In practice, we set  $\eta = 1$  for training (equivalent to applying DDPM [72]) and then  $\eta = 0$  for evaluation. We use DDIM sampling for our pixel-based experiments and long-horizon furniture assembly tasks, where the efficiency improvements are much desired.

**Diffusion noise scheduling.** We use the cosine schedule for  $\sigma_k$  introduced in [49], which was originally annealed to a small value on the order of  $1e-4$  at  $k = 0$ . In **DPO**, the values of  $\sigma_k$  also translate to the exploration noise that is crucial to training efficiency. Empirically, we find that clipping  $\sigma_k$  to a higher minimum value (denoted  $\sigma_{\min}^{\text{exp}}$ , e.g.,  $0.01 - 0.1$ ) when sampling actions helps exploration (see study in [Appendix B.1](#)). Additionally we clip  $\sigma_k$  to be at least 0.1 (denoted  $\sigma_{\min}^{\text{prob}}$ ) when evaluating the Gaussian likelihood  $\log \bar{\pi}_{\theta}(\bar{a}_{\bar{t}}|\bar{s}_{\bar{t}})$ , which improves training stability by avoiding large magnitude.

---

<sup>3</sup>In practice, we use Generalized Advantage Estimation (GAE) [69] that better balances variance and bias in estimating the advantage. We present the simpler form here.

**Network architecture.** We study both Multi-layer Perceptron (MLP) and UNet [66] as the policy heads in Diffusion Policy. MLP offers a simpler setup, and we find that it generally fine-tunes more stably with **Dppo**. Meanwhile, UNet, only applying convolutions to  $a_t^k$ , has the benefit of allowing pre-training and fine-tuning with different action chunk size  $T_a$ , e.g., 16 and 8. We find that **Dppo** benefits from pre-training with larger  $T_a$  (better prediction) and fine-tuning with smaller  $T_a$  (more amenable to policy gradient)<sup>4</sup>.

## 5 Performance Evaluation of **Dppo**

We begin our investigation of **Dppo** by studying its performance in various popular RL and robotics benchmarking environments. We compare to other RL approaches for fine-tuning a Diffusion Policy (Section 5.1), next to other policy parameterizations (Section 5.2), and then in multi-stage manipulation tasks including hardware evaluation (Section 5.3), and conclude with ablations (Section 5.4). While our evaluations focus primarily on **fine-tuning**, we also present training-from-scratch results in Appendix B. Wall-clock times are reported and discussed in Appendix C; they are roughly comparable (and often faster) than other diffusion-based RL baselines, though can be up to 2 $\times$  slower than other policy parameterization baselines. Full choices of training hyperparameters and additional training details are presented in Appendix D.

**Environments: OpenAI Gym.** We first consider three OpenAI GYM locomotion benchmarks [8] common in the RL literature: {Hopper-v2, Walker2D-v2, HalfCheetah-v2}. All policies are pre-trained with the medium-level datasets from D4RL [24] with **state** input and action chunk size  $T_a = 4$ . We use the original **dense** reward setup in fine-tuning.

**Environments: Robomimic.** Next we consider four simulated robot manipulation tasks from the ROBOMIMIC benchmark [48], {Lift, Can, Square, Transport}, ordered in increasing difficulty. These tasks are more representative of real-world robotic tasks, and Square and Transport (Fig. 4) are considered very challenging for RL training. Both **state** and **pixel** policy input are considered. State-based and pixel-based policies are pre-trained with 300 and 100 Multi-Human demonstrations provided by ROBOMIMIC, respectively. We consider  $T_a = 4$  for Can, Lift, and Square, and  $T_a = 8$  for Transport. They are then fine-tuned with **sparse** reward equal to 1 upon task completion.

**Environments: Furniture-Bench & real furniture assembly.** Finally, we demonstrate solving longer-horizon, multi-stage robot manipulation tasks from the FURNITURE-BENCH [28] benchmark in both simulation and reality. We evaluate the methods on three simulated furniture assembly tasks, {One-leg, Lamp, Round-table}, shown in Fig. 4 and described in detail in Appendix D.7. We consider two levels of randomness over initial state distribution, Low and Med, defined by the benchmark. **State** policy input is considered. All policies are pre-trained with 50 human demonstrations collected in simulation and  $T_a = 8$ . They are then fine-tuned with **sparse** (indicator of task stage completion) reward. We also evaluate the **zero-shot sim-to-real performance** with One-leg.

### 5.1 Comparison to diffusion-based RL algorithms

We compare **Dppo** to an extensive list of RL methods for fine-tuning diffusion models. Baseline names are color-coordinated with their plot colors. Two methods, **DRWR**, **DAWR**, are *our own, novel* baselines that are based on reward-weighted regression (RWR) [55] and advantage-weighted regression (AWR) [53]. The remaining methods, **DIPO** [86], **IDQL** [27], **DQL** [84], and **QSM** [59], are existing in the literature.

---

<sup>4</sup>With fully-connected layers in MLP, empirically we find that using different chunk sizes for pre-training and fine-tuning with MLP leads to training instability.

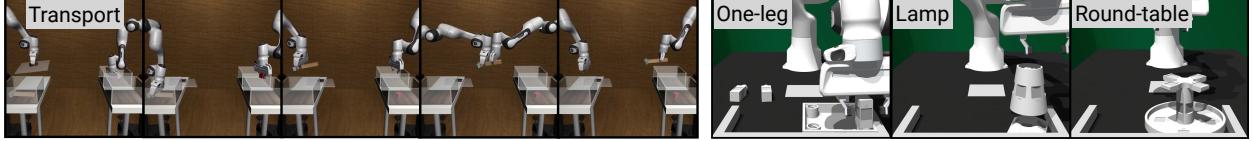


Figure 4: **Long-horizon robot manipulations tasks** including (left) the bimodal Transport from ROBOMIMIC and (right) FURNITURE-BENCH tasks (full rollouts visualized in Fig. A10).

We evaluate on the three OpenAI GYM tasks and the four ROBOMIMIC tasks with **state** input; detailed descriptions of all baselines and training details are in Appendix D.3.

Overall, **Dppo** performs consistently, exhibits great training stability, and enjoys strong fine-tuning success rates across tasks. In the GYM tasks (Figure 5, top row), **IDQL** and **DAWR** exhibit competitive performance, while the other methods perform worse and train less stably. **Dppo** is the strongest performer in the ROBOMIMIC tasks (Figure 5, bottom row), especially in the challenging Transport tasks. Surprisingly, **DRWR** is a strong baseline in {Lift, Can, Square} but underperforms in Transport, while all other baselines fare worse still. We postulate that the other baselines, all performing off-policy Q-function-based updates, can suffer from training instability in sparse-reward ROBOMIMIC tasks especially with large action chunk sizes (see further studies in Appendix B.1).

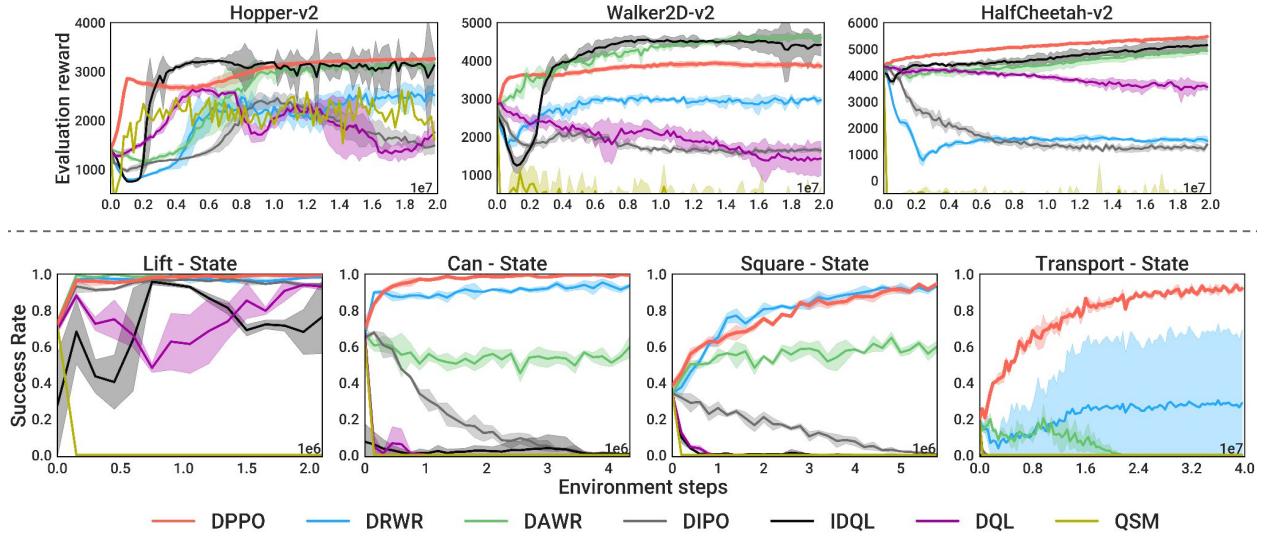


Figure 5: **Comparing to other diffusion-based RL algorithms.** Top row: GYM tasks [8]. Bottom row: ROBOMIMIC tasks [48] with **state** observation. Results are averaged over five seeds in GYM tasks and three seeds in ROBOMIMIC tasks. We opt to not show the erratic error bar with **QSM** in Hopper-v2. **Dppo** curves are slightly thicker for better visualization.

## 5.2 Comparison to other policy parameterizations

We compare **Dppo** with popular RL policy parameterizations: unimodal Gaussian with diagonal covariance [76] and Gaussian Mixture Model (GMM) [5], using either MLPs or Transformers [79], and also fine-tuned with the PPO objective. We compare these to **Dppo**-MLP and **Dppo**-UNet, which use either MLP or UNet as the network backbone. We evaluate on the four tasks from Robomimic (Lift, Can, Square, Transport) with both **state** and **pixel** input. With state input, **Dppo** pre-trains with 20 denoising steps

and then fine-tunes the last 10. With pixel input, **Dppo** pre-trains with 100 denoising steps and then fine-tunes 5 DDIM steps.

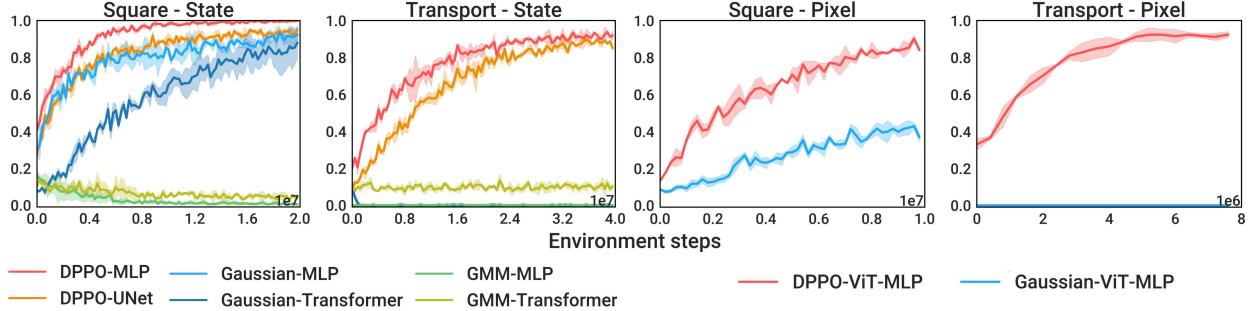


Figure 6: **Comparing to other policy parameterizations** in the more challenging Square and Transport tasks from ROBOMIC, with **state** (left) or **pixel** (right) observation. Results are averaged over three seeds.

**Fig. 6** shows the results in the more challenging Square and Transport — we defer the results in Lift and Can to **Fig. A8**. With **state** input, we find that **Dppo** outperforms Gaussian and GMM policies, with faster convergence to  $\sim 100\%$  success rate in Lift and Can, and greater final performance on Square and the challenging Transport, where it reaches  $> 90\%$ . UNet and MLP variants perform similarly, with the latter training somewhat more rapidly. With **pixel** inputs, we use a Vision-Transformer-based (ViT) image encoder introduced in Hu et al. [31] and an MLP head and compare the resulting variants **Dppo**-ViT-MLP and Gaussian-ViT-MLP (we omit GMM due to poor performance in state-based training). While the two are comparable on Lift and Can, **Dppo** trains more quickly and to higher accuracy on Square, and *drastically outperforms* on Transport, whereas Gaussian does not improve from its 0% pre-trained success rate. **To our knowledge, Dppo is the first RL algorithm to solve Transport from either state or pixel input to high ( $> 50\%$ ) success rates.**

### 5.3 Evaluation on Furniture-Bench, and sim-to-real transfer

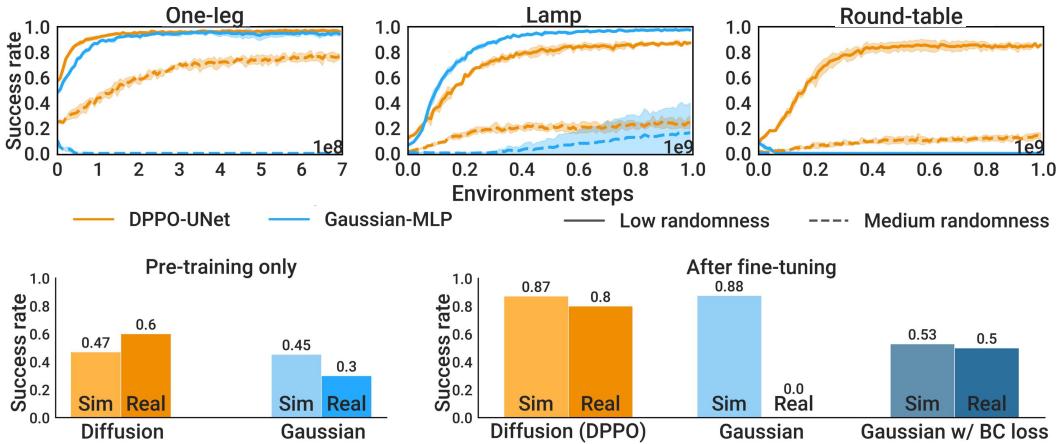


Figure 7: (Top) **Dppo** vs. Gaussian-MLP baseline in **simulated FURNITURE-BENCH tasks**. Results are averaged over three seeds. (Bottom) **Sim-to-real transfer results** in One-leg.

Here we evaluate **Dppo** on the long-horizon manipulation tasks from FURNITURE-BENCH [28]. We

compare **DPO** to Gaussian-MLP, the overall most effective baseline from [Section 5.2](#). [Fig. 7](#) (top row) shows the evaluation success rate over fine-tuning iterations. **DPO** exhibits strong training stability and improves policy performance in all six settings. Gaussian-MLP collapses to zero success rate in all three tasks with Med randomness (except for one seed in Lamp) and Round-table with Low randomness.

Note that we are only using 50 human demonstrations for pre-training; we expect **DPO** can leverage additional human data (better state space coverage) to further improve in Med, which is corroborated by ablation studies in [Appendix B.2](#) examining the effect of pre-training data on **DPO**.



**Figure 8: Qualitative comparison of pre-trained vs. fine-tuned DPO policies in hardware evaluation.** **(A)** Successful rollout with the pre-trained policy. **(B)** Failed rollout with the pre-trained policy due to imprecise insertion. **(C)** Successful rollout with the fine-tuned policy. **(D)** Successful rollout with the fine-tuned policy exhibiting corrective behavior.

**Sim-to-real transfer.** We evaluate **DPO** and Gaussian policies trained in the simulated One-leg task on physical hardware zero-shot (i.e., **no real data fine-tuning / co-training**) over 20 trials. Please see additional simulation training and hardware details in [Appendix D.7](#). [Fig. 7](#) (bottom row) shows simulated and hardware success rates after pre-training and fine-tuning. Notably, **DPO** improves the real-world performance to 80% (16 out of 20 trials). Though the Gaussian policy achieves a high success rate in simulation after fine-tuning (88%), it fails entirely on hardware (0%). Supplemental video suggests it exhibits volatile and jittery behavior. For stronger comparison, we also fine-tune the Gaussian policy with an auxiliary behavior-cloning loss [78] such that the fine-tuned policy is encouraged to stay close to the base policy. However, this limits fine-tuning and only leads to a 53% success rate in simulation and 50% in reality.

Qualitatively, we find fine-tuned policies to be more robust and exhibit more corrective behaviors than pre-trained-only policies, especially during the insertion stage of the task. Fig. 8 shows frames from representative rollouts on hardware. Overall, these results demonstrate the strong sim-to-real capabilities of **DPO**; Section 6 provides a conjectural mechanism for why this may be the case.

## 5.4 Summary of ablation findings

We conduct extensive ablation studies in Appendix B.1. Our main findings include: (1) for challenging tasks, using a value estimator which depends on environment state but is *independent of denoised action* is crucial for performance; we conjecture that this is related to the high stochasticity of Diffusion Policy; (2) there is a sweet spot for clipping the denoising noise level for **DPO** exploration, trading off between too little exploration and too much action noise; (3) **DPO** is resilient to fine-tuning fewer-than- $K$  denoising steps, yielding improved runtime and comparable performance; (4) **DPO** yields improvements over Gaussian-MLP baselines for varying levels of expert demonstration data, and achieves comparable final performance and sample efficiency when **training from scratch** in GYM environments.

## 6 Understanding the performance of **DPO**

We study the factors contributing to **DPO**'s improvements in performance over the popular Gaussian and GMM methods introduced in Section 5.2. Our findings highlight three major contributing factors:

- (1) **DPO** induces structured exploration near the pre-training data manifold [21].
- (2) **DPO** updates the action distribution progressively through the multi-step denoising process, which can be flexible and robust to policy collapse.
- (3) **DPO** leads to fine-tuned policies robust to perturbations in dynamics and initial state distribution.

We use the `Avoid` environment from the D3IL benchmark introduced in Jia et al. [36], where a robot arm needs to reach the other side of the table while avoiding an array of obstacles (Fig. 9, top-left). The action space is the 2D target location of the end-effector. D3IL provides a set of expert demonstrations that covers different possible paths to the goal line — we consider three subsets of the demonstrations, M1, M2, and M3 in Fig. 9, each with two distinct modes. We choose such relatively simple setups with only two modes in each setting such that Gaussian (with exploration noise)<sup>5</sup> and GMM can fit the expert data distribution reasonably well, allowing fair comparisons in fine-tuning.

We pre-train MLP-based Diffusion, Gaussian, and GMM policies (action chunk size  $T_a = 4$  unless noted) with the demonstrations. For fine-tuning, we assign (sparse) reward when the robot reaches the goal line from the topmost mode. Gaussian and GMM policies are also fine-tuned with the PPO objective.

**Benefit 1: Structured, on-manifold exploration.** Fig. 9 (right) shows the sampled trajectories from **DPO**, Gaussian, and GMM at the first iteration of fine-tuning. **DPO** explores in wide coverage **around the expert data manifold**, whereas Gaussian generates less structured exploration noise (especially in M2) and GMM exhibits narrower coverage. Moreover, the combination of diffusion parameterization with the denoising of *action chunks* means that policy stochasticity in **DPO** is **structured in both action dimension and time horizon**. Quantitatively, Fig. A7 in Appendix shows **DPO** achieves greater fine-tuning efficiency than Gaussian and GMM in all settings if a sufficient number of denoising steps is fine-tuned, consistent with the findings in Section 5.2 and Section 5.4.

It is possible, however, that the on-manifold exploration we observe with **DPO** hinders fine-tuning when aggressive, unstructured exploration is desired. We conjecture this is the case in the `Lamp` environment

---

<sup>5</sup>Without noise, Gaussian policy is fully deterministic and cannot capture the two modes.

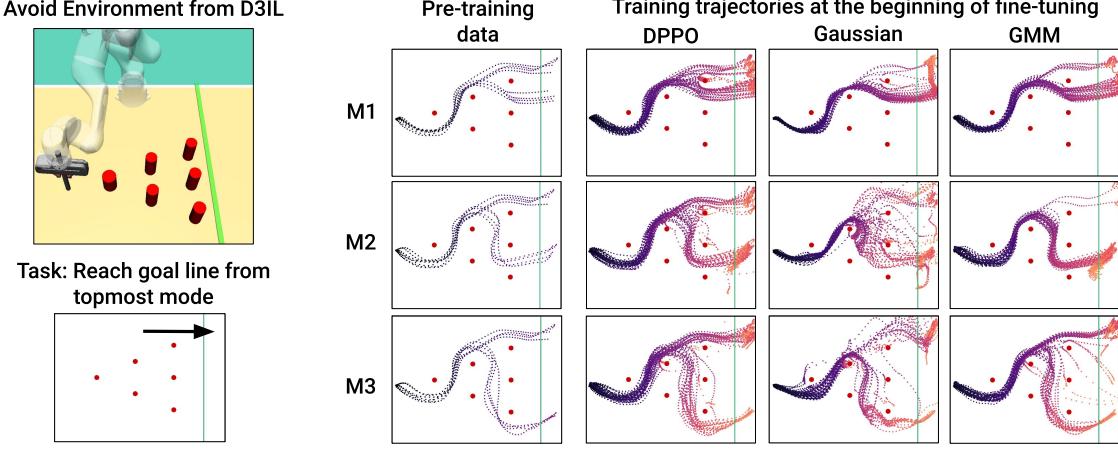


Figure 9: (Left) We use the **Avoid** environment from D3IL benchmark [36] to visualize the **DPPO**'s exploration tendencies. We design the task of always reaching the green goal line from the topmost mode. (Right) **Structured exploration.** We show the sampled trajectories at the beginning of fine-tuning for DPPO, Gaussian, and GMM after pre-training on three sets of expert demonstrations, M1, M2, and M3.

with `Low` randomness, in which **DPPO** slightly underperforms the Gaussian baseline (Fig. 7). Moreover, we do not observe that **DPPO** is substantially better (nor is it any worse) than Gaussian in exploration **from-scratch** (see Appendix B). Thus, we anticipate that this structured, on-manifold exploration confers the greatest benefit when pre-training provides sufficient coverage of relevant success modes. In particular, we believe that this makes **DPPO** uniquely suited for fine-tuning large Diffusion Policy pre-trained on multiple tasks [83], as it may exhibit sufficient mode coverage given training data diversity.

**Benefit 2: Training stability from multi-step denoising process.** In Fig. 10 (left), we run fine-tuning after pre-training with M2 and attempt to *de-stabilize* fine-tuning by gradually adding noise to the action during the fine-tuning process (see Appendix D.8 for details). We find that Gaussian and GMM's performance both collapse, while with **DPPO**, the performance is robust to the noise if at least four denoising steps are used. This property also allows **DPPO** to apply significant noise to the sampled actions, simulating an imperfect low-level controller to facilitate sim-to-real transfer in Section 5.3. In Fig. 10 (right), we also find **DPPO** enjoys greater training stability when fine-tuning long action chunks, e.g., up to  $T_a = 16$ , while Gaussian and GMM can fail to improve at all.

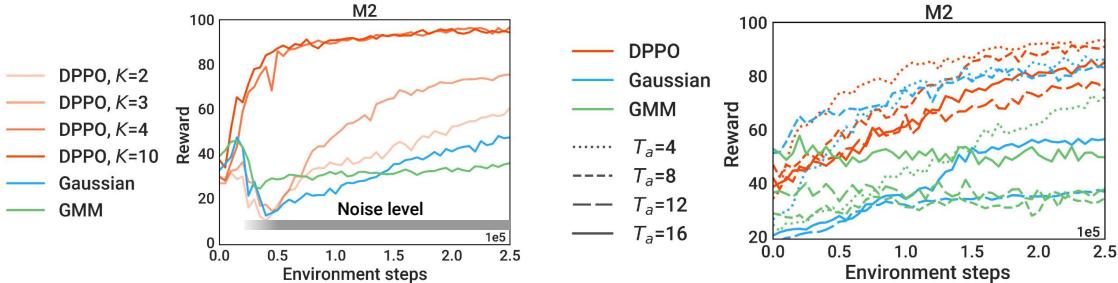


Figure 10: **Training stability.** Fine-tuning performance (averaged over five seeds, standard deviation not shown) after pre-training with M2. (Left) Noise is injected into the applied actions after a few training iterations. (Right) The action chunk size  $T_a$  is varied.

**Fig. 11** visualizes how **Dppo** affects the multi-step denoising process. Over fine-tuning iterations, the action distribution gradually converges through the denoising steps — the iterative refinement is largely preserved, as opposed to, e.g., “collapsing” to the optimal actions at the first fine-tuned denoising step or the final one. We postulate this contributes to the training stability of **Dppo**.

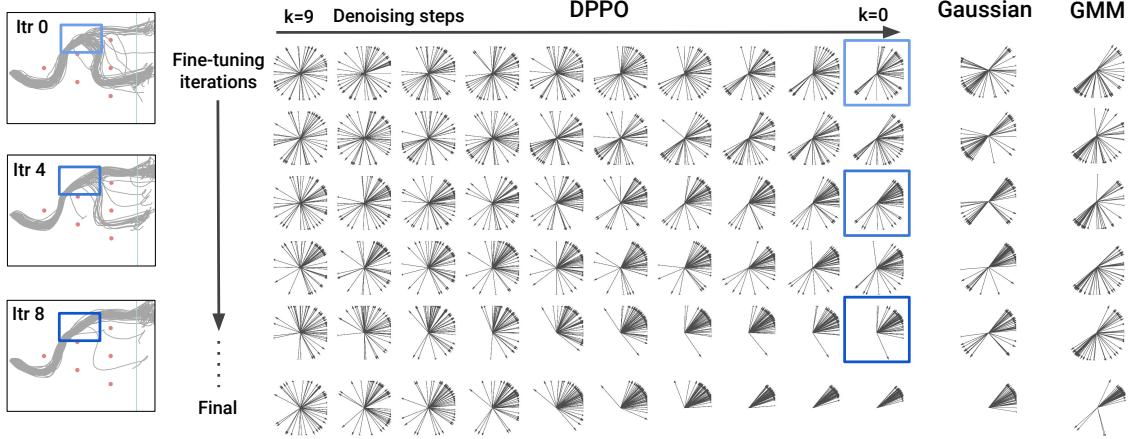


Figure 11: **Preserving the iterative refinement.** The 2D actions from 50 trajectories at the branching point through fine-tuning iterations after pre-training with M2. For **Dppo**, we also visualize the action distribution through the final denoising steps at each fine-tuning iteration.

**Benefit 3: Robust and generalizable fine-tuned policy.** **Dppo** also generates final policies robust to perturbations in dynamics and the initial state distribution. In Fig. 12, we again add noise to the actions sampled from the fine-tuned policy (no noise applied during training) and find that the **Dppo** policy exhibits strong robustness to the noise compared to the Gaussian policy. The **Dppo** policy also converges to the (near-)optimal path from a larger distribution of initial states. This finding echoes theoretical guarantees that Diffusion Policy, capable of representing complex multi-modal data distribution, can effectively deconvolve noise from noisy states [7], a property used in Chen et al. [11] to stabilize long-horizon video generation.

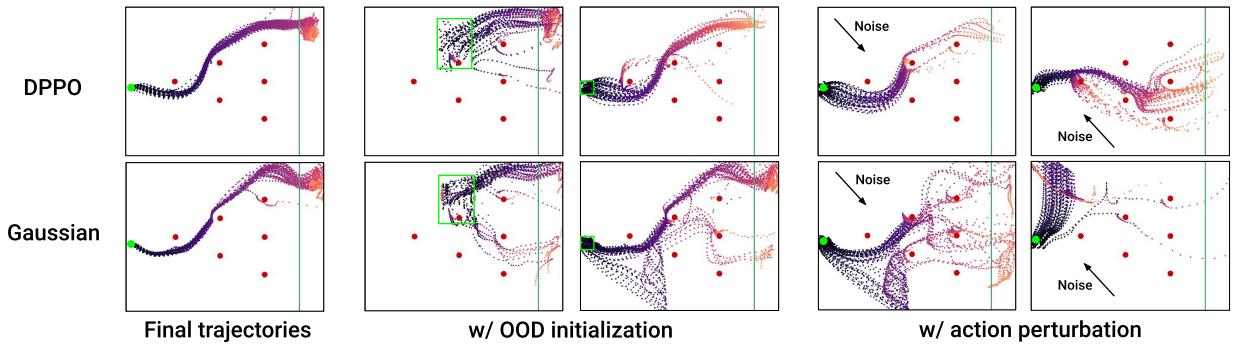


Figure 12: **Policy robustness.** Green dot / box indicates the initial state region.

## 7 Conclusion and Future Work

We present *Diffusion Policy Policy Optimization (Dppo)* for fine-tuning a pre-trained Diffusion Policy with the policy gradient method. **Dppo** leverages the sequential nature of the diffusion denoising process and fine-tunes the entire chain of diffusion MDPs. **Dppo** exhibits structured online exploration, strong training

stability, and robustness and generalization at deployment. We demonstrate the efficiency and effectiveness of **DPPO** fine-tuning in various RL and robotics benchmarks, as well as strong sim-to-real transfer of a **DPPO** policy in a long-horizon, multi-stage manipulation task.

We believe **DPPO** will become an important component in the pre-training-plus-fine-tuning pipeline for training general-purpose real-world robotic policies. To this end, we hope in future work to further showcase the promise of **DPPO** for simulation-to-real transfer [41, 62, 14] in which we fine-tune a vision-based policy that has been pre-trained on a variety of diverse tasks. We expect this pre-training to provide a large and diverse expert data manifold, of which, as we have shown in Section 6, **DPPO** is well-suited to take advantage of better exploration during fine-tuning. We are also excited to understand how **DPPO** can fit together with other decision-making tools, such as model-based planning [35] and decision-making aided by video prediction [11]. Finally, as noted in the introduction, we eagerly anticipate applications of **DPPO** in domains beyond robotics, where diffusion models have shown promise for combinatorial search and sequence modeling [58, 45].

## Acknowledgments

We would like to thank Lirui Wang and Terry Suh for helpful discussions in the early stage of the project. The authors were partially supported by the Toyota Research Institute (TRI). This article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

## References

- [1] J. Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- [2] A. Ajay, Y. Du, A. Gupta, J. B. Tenenbaum, T. S. Jaakkola, and P. Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh International Conference on Learning Representations*, 2023.
- [3] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 2020.
- [4] L. Ankile, A. Simeonov, I. Shenfeld, and P. Agrawal. Juicer: Data-efficient imitation learning for robotic assembly. *arXiv*, 2024.
- [5] C. M. Bishop and N. M. Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.
- [6] K. Black, M. Janner, Y. Du, I. Kostrikov, and S. Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- [7] A. Block, A. Jadbabaie, D. Pfommer, M. Simchowitz, and R. Tedrake. Provable guarantees for generative behavior cloning: Bridging low-level stability and high-level behavior. *Advances in Neural Information Processing Systems*, 2024.
- [8] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [9] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 2020.
- [10] S. H. Chan. Tutorial on diffusion models for imaging and vision. *arXiv preprint arXiv:2403.18103*, 2024.

- [11] B. Chen, D. M. Monso, Y. Du, M. Simchowitz, R. Tedrake, and V. Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *arXiv preprint arXiv:2407.01392*, 2024.
- [12] H. Chen, C. Lu, C. Ying, H. Su, and J. Zhu. Offline reinforcement learning via high-fidelity generative behavior modeling. *arXiv preprint arXiv:2209.14548*, 2022.
- [13] T. Chen, J. Xu, and P. Agrawal. A system for general in-hand object re-orientation. In *Conference on Robot Learning*, 2022.
- [14] C. Chi, B. Burchfiel, E. Cousineau, S. Feng, and S. Song. Iterative residual policy: for goal-conditioned dynamic manipulation of deformable objects. *The International Journal of Robotics Research*, 2024.
- [15] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 2024.
- [16] K. Clark, P. Vicol, K. Swersky, and D. J. Fleet. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- [17] Z. Ding and C. Jin. Consistency models as a rich and efficient policy class for reinforcement learning. *arXiv preprint arXiv:2309.16984*, 2023.
- [18] L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry. Implementation matters in deep rl: A case study on ppo and trpo. In *International conference on learning representations*, 2019.
- [19] Y. Fan and K. Lee. Optimizing ddpm sampling with shortcut fine-tuning. *arXiv preprint arXiv:2301.13362*, 2023.
- [20] Y. Fan, O. Watkins, Y. Du, H. Liu, M. Ryu, C. Boutilier, P. Abbeel, M. Ghavamzadeh, K. Lee, and K. Lee. Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 2024.
- [21] C. Fefferman, S. Mitter, and H. Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 2016.
- [22] P. Florence, L. Manuelli, and R. Tedrake. Self-supervised correspondence in visuomotor policy learning. *IEEE Robotics and Automation Letters*, 2019.
- [23] P. Florence, C. Lynch, A. Zeng, O. A. Ramirez, A. Wahid, L. Downs, A. Wong, J. Lee, I. Mordatch, and J. Tompson. Implicit behavioral cloning. In *Conference on Robot Learning*. PMLR, 2022.
- [24] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- [25] Z. Fu, T. Z. Zhao, and C. Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024.
- [26] W. Goo and S. Niekum. Know your boundaries: The necessity of explicit behavioral cloning in offline rl. *arXiv preprint arXiv:2206.00695*, 2022.
- [27] P. Hansen-Estruch, I. Kostrikov, M. Janner, J. G. Kuba, and S. Levine. Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint arXiv:2304.10573*, 2023.
- [28] M. Heo, Y. Lee, D. Lee, and J. J. Lim. Furniturebench: Reproducible real-world benchmark for long-horizon complex manipulation. *arXiv preprint arXiv:2305.12821*, 2023.
- [29] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 2020.
- [30] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- [31] H. Hu, S. Mirchandani, and D. Sadigh. Imitation bootstrapped reinforcement learning. *arXiv preprint arXiv:2311.02198*, 2023.

- [32] S. Huang, R. F. J. Dossa, C. Ye, J. Braga, D. Chakraborty, K. Mehta, and J. G. Araújo. Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 2022.
- [33] Z. Huang, L. Yang, X. Zhou, Z. Zhang, W. Zhang, X. Zheng, J. Chen, Y. Wang, C. Bin, and W. Yang. Protein-ligand interaction prior for binding-aware 3d molecule diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [34] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 2019.
- [35] M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior synthesis. *arXiv preprint arXiv:2205.09991*, 2022.
- [36] X. Jia, D. Blessing, X. Jiang, M. Reuss, A. Donat, R. Lioutikov, and G. Neumann. Towards diverse behaviors: A benchmark for imitation learning with human demonstrations. *arXiv preprint arXiv:2402.14606*, 2024.
- [37] B. Kang, X. Ma, C. Du, T. Pang, and S. Yan. Efficient diffusion policies for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 2024.
- [38] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza. Champion-level drone racing using deep reinforcement learning. *Nature*, 2023.
- [39] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- [40] S. Lee, Y. Wang, H. Etukuru, H. J. Kim, N. M. M. Shafiullah, and L. Pinto. Behavior generation with latent actions. *arXiv preprint arXiv:2403.03181*, 2024.
- [41] J. Liang, S. Saxena, and O. Kroemer. Learning active task-oriented exploration policies for bridging the sim-to-real gap. *arXiv preprint arXiv:2006.01952*, 2020.
- [42] Z. Liang, Y. Mu, M. Ding, F. Ni, M. Tomizuka, and P. Luo. Adaptdiffuser: Diffusion models as adaptive self-evolving planners. *arXiv preprint arXiv:2302.01877*, 2023.
- [43] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [44] Y. Lin, A. S. Wang, G. Sutanto, A. Rai, and F. Meier. Polymetis. <https://facebookresearch.github.io/fairo/polymetis/>, 2021.
- [45] A. Lou, C. Meng, and S. Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. *stat*, 2024.
- [46] S. Luo, Y. Su, X. Peng, S. Wang, J. Peng, and J. Ma. Antigen-specific antibody design and optimization with diffusion-based generative models for protein structures. *Advances in Neural Information Processing Systems*, 2022.
- [47] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- [48] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *arXiv preprint arXiv:2108.03298*, 2021.
- [49] A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, 2021.
- [50] T. Osa, J. Pajarinen, G. Neumann, J. A. Bagnell, P. Abbeel, J. Peters, et al. An algorithmic perspective on imitation learning. *Foundations and Trends® in Robotics*, 2018.

- [51] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 2022.
- [52] T. Pearce, T. Rashid, A. Kanervisto, D. Bignell, M. Sun, R. Georgescu, S. V. Macua, S. Z. Tan, I. Momennejad, K. Hofmann, et al. Imitating human behaviour with diffusion models. *arXiv preprint arXiv:2301.10677*, 2023.
- [53] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- [54] X. B. Peng, Z. Ma, P. Abbeel, S. Levine, and A. Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (ToG)*, 2021.
- [55] J. Peters and S. Schaal. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, 2007.
- [56] D. A. Pomerleau. Alvinn: An autonomous land vehicle in a neural network. *Advances in neural information processing systems*, 1988.
- [57] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [58] M. Popova, O. Isayev, and A. Tropsha. Deep reinforcement learning for de novo drug design. *Science advances*, 2018.
- [59] M. Psenka, A. Escontrela, P. Abbeel, and Y. Ma. Learning a diffusion model policy from rewards via q-score matching. *arXiv preprint arXiv:2312.11752*, 2023.
- [60] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 2021.
- [61] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- [62] A. Z. Ren, H. Dai, B. Burchfiel, and A. Majumdar. Adaptsim: Task-driven simulation adaptation for sim-to-real transfer. In *Proceedings of the Conference on Robot Learning (CoRL)*, 2023.
- [63] M. Reuss, M. Li, X. Jia, and R. Lioutikov. Goal-conditioned imitation learning using score-based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- [64] M. Rigter, J. Yamada, and I. Posner. World models via policy-guided trajectory diffusion. *arXiv preprint arXiv:2312.08533*, 2023.
- [65] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022.
- [66] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention (MICCAI)*, 2015.
- [67] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023.
- [68] S. S. Sahoo, M. Arriola, Y. Schiff, A. Gokaslan, E. Marroquin, J. T. Chiu, A. Rush, and V. Kuleshov. Simple and effective masked diffusion language models. *arXiv preprint arXiv:2406.07524*, 2024.
- [69] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.

- [70] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [71] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, 2015.
- [72] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [73] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [74] A. Sridhar, D. Shah, C. Glossop, and S. Levine. Nomad: Goal masked diffusion policies for navigation and exploration. *arXiv preprint arXiv:2310.07896*, 2023.
- [75] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [76] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 1999.
- [77] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, 2012.
- [78] M. Torne, A. Simeonov, Z. Li, A. Chan, T. Chen, A. Gupta, and P. Agrawal. Reconciling reality through simulation: A real-to-sim-to-real approach for robust manipulation. *arXiv preprint arXiv:2403.03949*, 2024.
- [79] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 2017.
- [80] S. Venkatraman, S. Khaitan, R. T. Akella, J. Dolan, J. Schneider, and G. Berseth. Reasoning with latent diffusion in offline reinforcement learning. *arXiv preprint arXiv:2309.06599*, 2023.
- [81] B. Wallace, M. Dang, R. Rafailov, L. Zhou, A. Lou, S. Purushwalkam, S. Ermon, C. Xiong, S. Joty, and N. Naik. Diffusion model alignment using direct preference optimization. *arXiv preprint arXiv:2311.12908*, 2023.
- [82] J. Wang and E. Olson. Apriltag 2: Efficient and robust fiducial detection. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016.
- [83] L. Wang, J. Zhao, Y. Du, E. H. Adelson, and R. Tedrake. Poco: Policy composition from and for heterogeneous robot learning. *arXiv preprint arXiv:2402.02511*, 2024.
- [84] Z. Wang, J. J. Hunt, and M. Zhou. Diffusion policies as an expressive policy class for offline reinforcement learning. *arXiv preprint arXiv:2208.06193*, 2022.
- [85] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992.
- [86] L. Yang, Z. Huang, F. Lei, Y. Zhong, Y. Yang, C. Fang, S. Wen, B. Zhou, and Z. Lin. Policy representation via diffusion probability model for reinforcement learning. *arXiv preprint arXiv:2305.13122*, 2023.
- [87] Y. Ze, G. Zhang, K. Zhang, C. Hu, M. Wang, and H. Xu. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. In *ICRA 2024 Workshop on 3D Visual Representations for Robot Manipulation*.
- [88] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- [89] Z. Zhu, H. Zhao, H. He, Y. Zhong, S. Zhang, Y. Yu, and W. Zhang. Diffusion models for reinforcement learning: A survey. *arXiv preprint arXiv:2311.01223*, 2023.

## A Extended Related Work

This section presents an extended discussion on related methods that train/improve diffusion-based policies with RL-related methods. The baselines to which we compare in [Section 5.1](#) are discussed below as well, and are highlighted in their corresponding colors. We also refer the readers to Zhu et al. [89] for an extensive survey on diffusion models for RL.

Most previous work have focused on the **offline** setting with a static dataset. One line of work focuses on state trajectory planning and *guiding* the denoising sampling process such that the sampled actions satisfy some desired objectives. Janner et al. [35] applies classifier guidance that generates trajectories with higher predicted reward. Ajay et al. [2] introduces classifier-free guidance that avoids learning the value of noisy states. There is another line of work that uses diffusion model as action policy (instead of state planner) and generally applies Q-learning. **DQL** [84] introduces Diffusion Q-Learning that learns a state-action critic for the final denoised actions and backpropagates the gradient from the critic through the entire Diffusion Policy (actor), akin to the usual Q-learning. **IDQL** [27], or Implicit Diffusion Q-learning, proposes learning the critic to select the actions at inference time for either training or evaluation, while fitting the actor to all sampled actions. Kang et al. [37] instead proposes using the critic to re-weight the sampled actions for updating the actor itself, similar to weighted regression baselines **DAWR** and **DRWR** introduced in our work. Goo and Niekum [26] similarly extracts the policy in the spirit of AWR [53]. Chen et al. [12] trains the critic using value iteration instead based on samples from the actor.

We note that methods like **DQL** and **IDQL** can also be applied in the **online** setting. There is also a small amount of work that focuses entirely on the online setting. **DIPPO** [86] differs from **DQL** and related work that it uses the critic to update the sampled actions (“action gradient”), instead of the actor — the actor is then fitted with updated actions from the replay buffer. **QSM**, or Q-Score Matching [59], suggests that optimizing the likelihood of the entire chain of denoised actions can be inefficient (contrary to our findings in the fine-tuning setting), and instead proposes learning the optimal policy by iteratively aligning the gradient of the actor (i.e., score) with the action gradient of the critic. Rigter et al. [64] proposes learning a diffusion dynamic model to generate synthetic trajectories for training a non-diffusion RL policy online.

## B Additional experimental results

### B.1 Ablation studies on design decisions in **DPPO**

**1. Choice of advantage estimator.** In [Section 4.3](#) we demonstrate how to efficiently estimate the advantage used in PPO updates by learning  $\tilde{V}(s_t)$  that only depends on the environment state; the advantage used in **DPPO** is formally

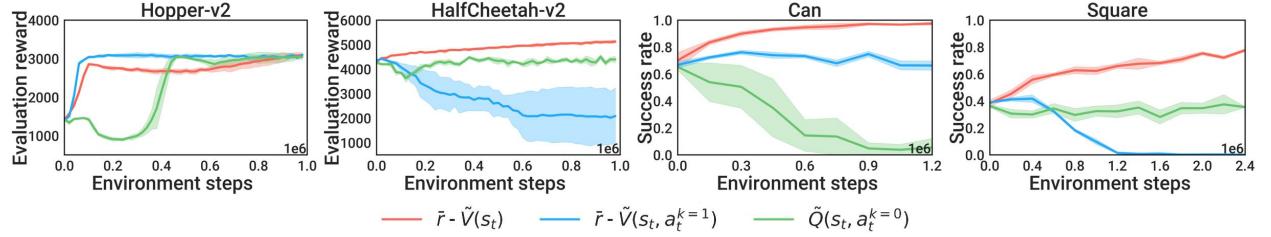
$$\hat{A} = \gamma_{\text{DENOISE}}^k(\bar{r}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) - \tilde{V}(s_t)).$$

We now compare this choice with learning the value of the full state  $\bar{s}_{\bar{t}(t,0)}$  that includes environment state  $s_t$  and denoised action  $a_t^{k=1}$ . We additionally compare with the state-action Q-function estimator used in Psenka et al. [59]<sup>6</sup>,  $\tilde{Q}(s_t, a_t^{k=0})$ , that does not directly use the rollout reward  $\bar{r}$  in the advantage.

[Fig. A1](#) shows the fine-tuning results in Hopper-v2 and HalfCheetah-v2 from GYM, and Can and Square from ROBOMIMIC. On the simpler Hopper-v2, we observe that the two baselines, both estimating the value of some action, achieves higher reward during fine-tuning than **DPPO**’s choice. However, in the more challenging tasks, the environment-state-only advantage used in **DPPO** consistently leads to the

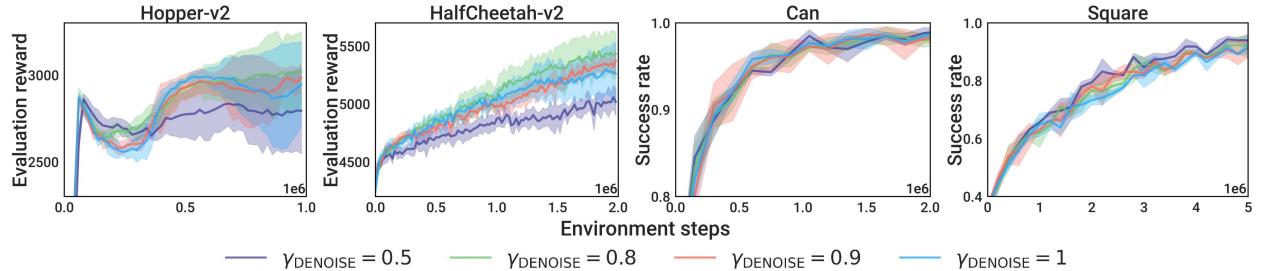
<sup>6</sup>Psenka et al. [59] applies off-policy training with double Q-learning (according to its open-source implementation) and policy gradient over the denoising steps. Note that this is a baseline in Psenka et al. [59] that is conjectured to be inefficient. We follow the same except for applying on-policy PPO updates.

most improved performance. We believe estimating the accurate value of applying a continuous and high-dimensional action can be challenging, and this is exacerbated by the high stochasticity of diffusion-based policies and the action chunk size. The results here corroborate the findings in [Section 5.1](#) that off-policy Q-learning methods can perform well in Hopper-v2 and Walker2D-v2, but exhibit training instability in manipulation tasks from ROBOMIMIC.



**Figure A1: Choice of advantage estimator.** Results are averaged over five seeds in Hopper-v2 and HalfCheetah-v2 and three seeds in Can and Square.

**Denoising discount factor.** We further examine how  $\gamma_{\text{DENOISE}}$  in the **DPPO** advantage estimator affects fine-tuning. Using a smaller value (i.e., more discount) has the effect of downweighting the contribution of earlier denoising steps in the policy gradient. [Fig. A2](#) shows the fine-tuning results in the same four tasks with varying  $\gamma_{\text{DENOISE}} \in [0.5, 0.8, 0.9, 1]$ . We find in Hopper-v2 and HalfCheetah-v2  $\gamma_{\text{DENOISE}} = 0.8$  leads to better efficiency while smaller  $\gamma_{\text{DENOISE}} = 0.5$  slows training. The value does not affect training noticeably in Can. In Square the smaller  $\gamma_{\text{DENOISE}} = 0.5$  works slightly better. Overall in manipulation tasks, **DPPO** training seems relatively robust to this choice.



**Figure A2: Choice of denoising discount factor.** Results are averaged over five seeds in Hopper-v2 and HalfCheetah-v2 and three seeds in Can and Square.

**2. Choice of diffusion noise schedule.** As introduced in [Section 4.3](#), we find it helpful to clip the diffusion noise  $\sigma_k$  to a higher minimum value  $\sigma_{\min}^{\text{exp}}$  to ensure sufficient exploration. In [Figure A3](#), we perform analysis on varying  $\sigma_{\min}^{\text{exp}} \in \{.001, .01, .1, .2\}$  (keeping  $\sigma_{\min}^{\text{prob}} = .1$  to evaluate likelihoods). Although in Can the choice of  $\sigma_{\min}^{\text{exp}}$  does not affect the fine-tuning performance, in Square a higher  $\sigma_{\min}^{\text{exp}} = 0.1$  is required to prevent the policy from collapsing. We conjecture that this is due to limited exploration causing policy over-optimizing the collected samples that exhibit limited state-action coverage. We also visualize the trajectories at the beginning of fine-tuning in Avoid task from D3IL. With higher  $\sigma_{\min}^{\text{exp}}$ , the trajectories still remain near the two modes of the pre-training data but exhibit a higher coverage in the state space — we believe this additional coverage leads to better exploration. Anecdotally, we find terminating the denoising process early can also provide exploration noise and lead to comparable results, but it requires a more involved implementation around the denoising MDP.

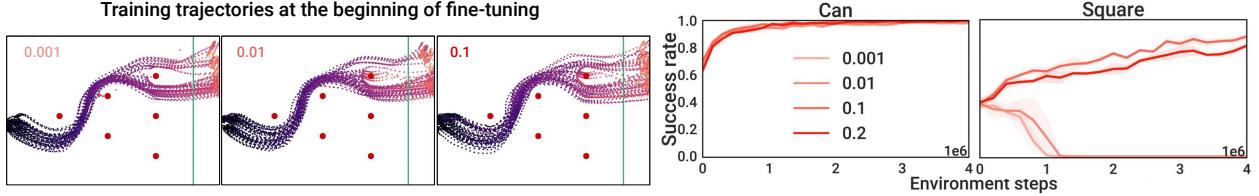


Figure A3: **Choice of minimum diffusion noise.** Results are averaged over three seeds.

**3. Choice of the number of fine-tuned denoising steps.** We examine how the number of fine-tuned denoising steps in **DPPo**,  $K'$ , affects the fine-tune performance and wall-clock time in Fig. A4. We show the curves of individual runs (three for each  $K'$ ) instead of the average as their wall-clock times (X-axis) are not perfectly aligned. Generally, fine-tuning too few denoising steps (e.g., 3) can lead to subpar asymptotic performance and slower convergence especially in Can. Fine-tuning 10 steps leads to the overall best efficiency. Similar results are also shown in Fig. A7 with Avoid task. Lastly, we note that the GPU memory usage scales linearly with  $K'$ .

We note that the findings here mostly correlate with those from varying the denoising discount factor,  $\gamma_{\text{DENOISE}}$ . Discounting the earlier denoising steps in the policy gradient can be considered as a soft version of hard limiting the number of fine-tuned denoising steps. Depending on the amount of fine-tuning needed from the pre-trained action distribution, one can flexibly adjust  $\gamma_{\text{DENOISE}}$  and  $K'$  to achieve the best efficiency.

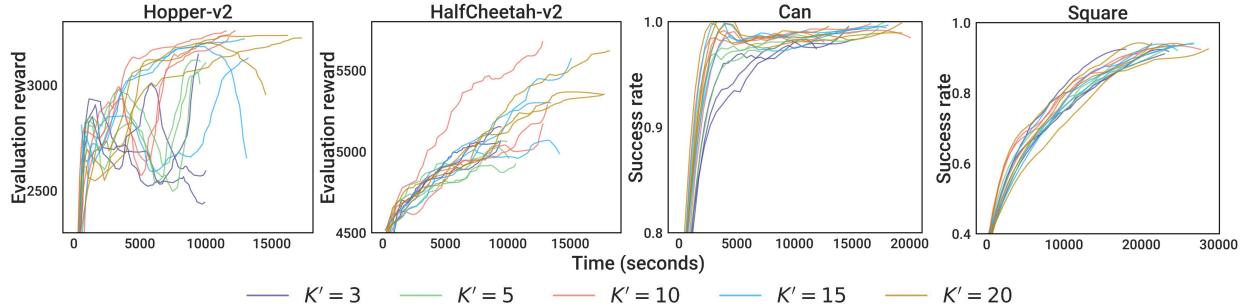


Figure A4: **Choice of number of fine-tuned denoising steps,  $K'$ .** Individual runs are shown. The curves are smoothed using a Savitzky–Golay filter.

## B.2 Effect of expert data

We investigate the effect of the amount of pre-training expert data on fine-tuning performance. In Fig. A5 we compare **DPPo** and Gaussian in Hopper-v2, Square, and One-leg task from FURNITURE-BENCH, using varying numbers of expert data (episodes) denoted in the figure. Overall, we find **DPPo** can better leverage the pre-training data and fine-tune to high success rates. Notably, **DPPo** obtains non-trivial performance (60% success rate) on One-leg from only 10 episode of demonstrations.

**Training from scratch.** In Fig. A6 we compare **DPPo** (10 denoising steps) and Gaussian *trained from scratch* (no pre-training on expert data) in the three OpenAI GYM tasks. As using larger action chunk sizes  $T_a$  leads to poor from-scratch training shown in Fig. A5, we focus on single-action chunks  $T_a = 1$  as is typical in RL benchmarking. Though we find Gaussian trains faster than **DPPo** (expected since **DPPo** solves an MDP with longer effective horizon), **DPPo** still attains reasonable final performance. However, due to the multi-step (10) denoising sampling, **DPPo** takes about  $6\times$  wall-clock time compared to Gaussian. We hope that future work will explore how to design the training curriculum of denoising steps for the best balance of training performance and wall-clock efficiency.

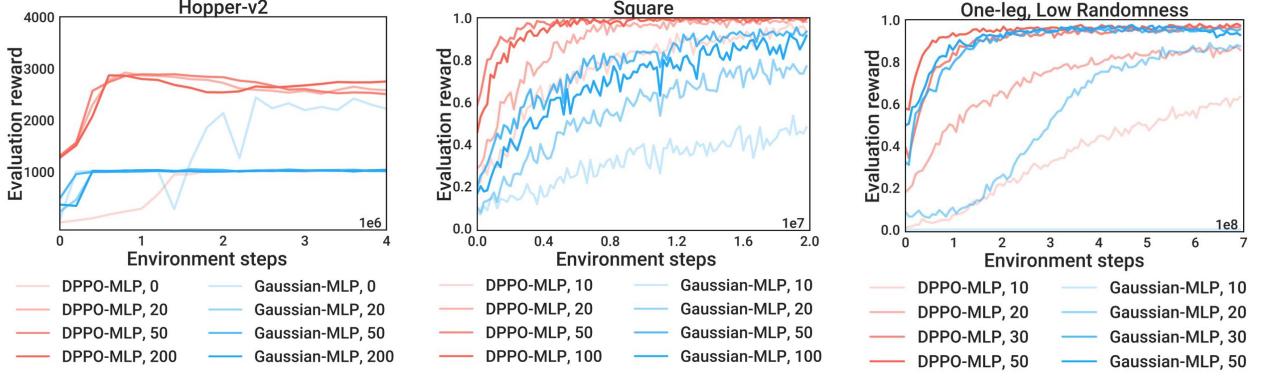


Figure A5: **Varying the number of expert demonstrations.** The numbers in the legends indicates the number of episodes used in pre-training.

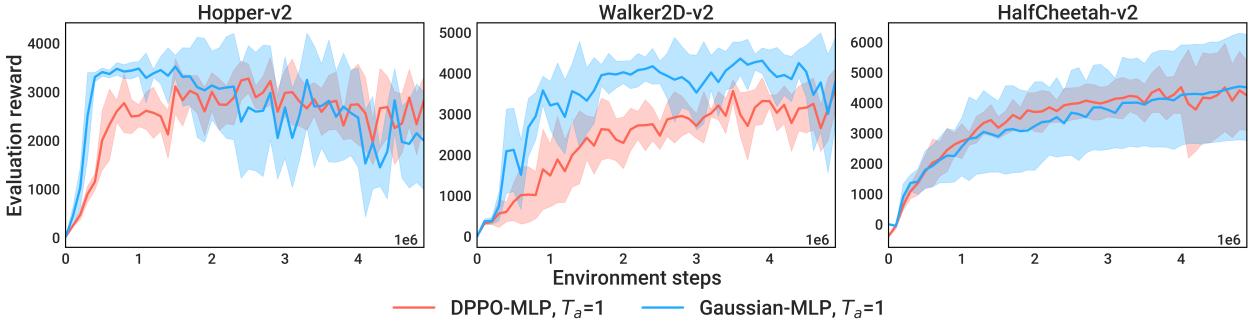


Figure A6: **No expert data / pre-training** with GYM tasks. Results are averaged over five seeds.

### B.3 Comparing to other policy parameterizations in **Avoid**

Figure A7 depicts the performance of various parameterizations of **Dppo** (with differing numbers of fine-tuned denoising steps,  $K'$ ) to Gaussian and GMM baselines. We study the **Avoid** task from D3IL, after pre-training with the data from M1, M2, M3 as described in Section 6. We find that, for  $K' \in \{15, 20\}$ , **Dppo** attains the highest performance of all methods and trains the quickest in terms of environment steps; on M1, M2, it appears to attain the greatest terminal performance as well.  $K' = 10$  appears slightly better than, but roughly comparable to, the Gaussian baseline, with GMM and  $K' < 10$  performing less strongly.

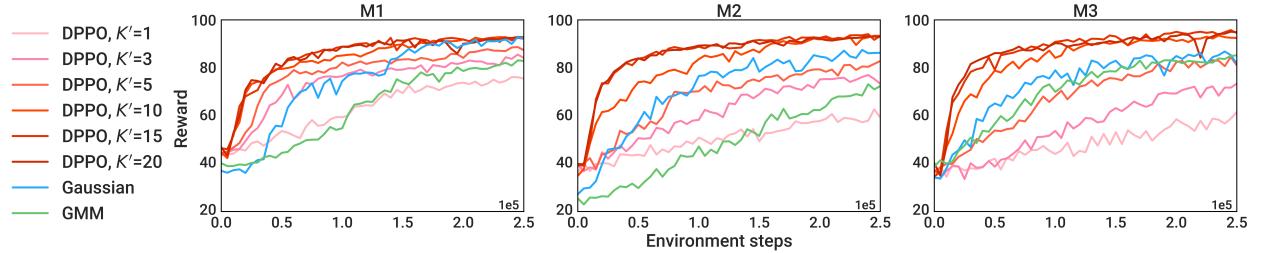


Figure A7: Fine-tuning performance (averaged over five seeds, standard deviation not shown) after pre-training with M1, M2, and M3 in **Avoid task from D3IL**. **Dppo** ( $K = 20$ ), Gaussian, and GMM policies are compared. We also sweep the number of fine-tuned denoising steps  $K'$  in **Dppo**.

## B.4 Comparing to other policy parameterizations in the easier tasks from ROBOMIMIC

Figure A8 compares the performance of **DPO** to Gaussian and GMM baselines, across a variety of architectures, and with **state** and **pixel** inputs, in Lift and Can environments in the ROBOMIMIC suite. Compared to the Square and Transport (results shown in Section 5), these environments are considered to be “easier”, and this is reflected in the greater performance of **DPO** and Gaussian baselines (GMM still exhibits subpar performance). Nonetheless, **DPO** still achieves similar or even better sample efficiency compared to Gaussian baseline.

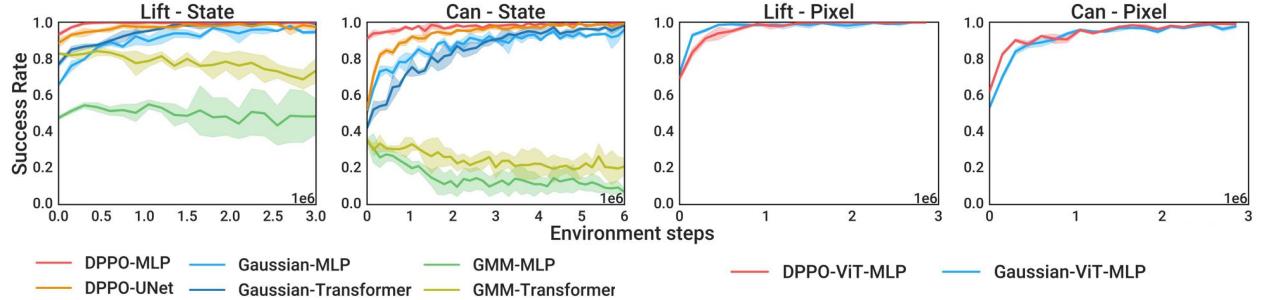


Figure A8: **Comparing to other policy parameterizations** in the easier Lift and Can tasks from ROBOMIMIC, with **state** (left) or **pixel** (right) observation. Results are averaged over three seeds.

## B.5 Comparing to policy gradient using exact likelihood of Diffusion Policy

Here we experiment another novel method (which, to our knowledge, has not been explicitly studied in any previous work) for performing policy gradient with diffusion-based policies. Although diffusion model does not directly model the action likelihood,  $p_\theta(a_0|s)$ , there have been ways to *estimate* the value, e.g., by solving the probability flow ODE that implements DDPM [73]. We refer the readers to Appendix. D in Song et al. [73] for a comprehensive exposition. We follow the official open-source code from Song et al.<sup>7</sup>, and implement policy gradient (single-level MDP) that uses the exact action likelihood  $\pi_\theta(a_t|s_t)$  (3.1).

Fig. A9 shows the comparison between **DPO** and diffusion policy gradient using exact likelihood estimate. Exact policy gradient improves the base policy in Hopper-v2 but does not outperform **DPO**. It also requires more runtime and GPU memory as it backpropagates through the ODE. In the more challenging Can its success rate drops to zero. Moreover, policy gradient with exact likelihood does not offer the flexibility of fine-tuning fewer-than- $K$  denoising steps or discounting the early denoising steps that **DPO** offers, which have shown in Appendix B.1 to often improve fine-tuning efficiency.

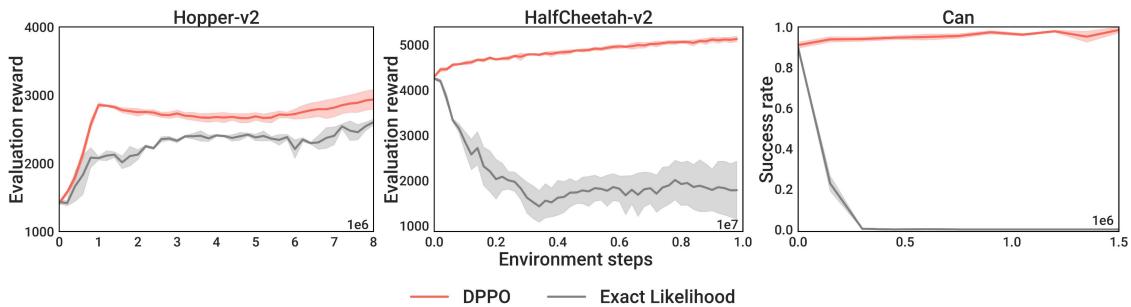


Figure A9: **Comparing to diffusion policy gradient with exact action likelihood**. Results are averaged over five seeds in Hopper-v2 and HalfCheetah-v2, and three seeds in Can.

<sup>7</sup>[https://github.com/yang-song/score\\_sde\\_pytorch](https://github.com/yang-song/score_sde_pytorch)

## C Reporting of Wall-Clock Times

**Comparing to other diffusion-based RL algorithms** **Section 5.1.** Table A1 and Table A2 shows the wall-clock time used in each OpenAI GYM task and ROBOMIMIC task. In GYM tasks, on average **DPPo** trains 74%, 41%, 37%, and 12% faster than **QSM**, **DAWR**, **DIPO**, and **DQL**, respectively, which all require a significant amount of gradient updates per sample to train stably. **DRWR** and **IDQL** trains 33% and 7% faster than **DPPo**, respectively. ROBOMIMIC tasks are more expensive to simulate, especially with Transport task, and thus the wall-clock difference is smaller among the different methods. On average **DPPo** trains 29%, 3%, 7%, and 3% faster than **QSM**, **DAWR**, **DIPO**, and **DQL**, respectively. **DRWR** trains 3% faster than **DPPo**, and **IDQL** uses comparable time as **DPPo**.

Method	Task		
	Hopper-v2	Walker2D-v2	HalfCheetah-v2
<b>DRWR</b>	11.3	12.7	10.4
<b>DAWR</b>	30.4	30.7	27.1
<b>DIPO</b>	27.8	27.9	26.0
<b>IDQL</b>	16.3	16.1	15.5
<b>DQL</b>	20.5	20.5	17.6
<b>QSM</b>	64.4	65.7	68.3
<b>DPPo</b>	16.6	18.3	16.8

Table A1: **Wall-clock time** in seconds for a single training iteration in **OpenAI GYM tasks** when comparing diffusion-based RL algorithms. Each iteration involves 500 environment timesteps in each of the 40 parallelized environments running on 40 CPU threads and a NVIDIA RTX 2080 GPU (20000 steps total).

Method	Task			
	Lift	Can	Square	Transport
<b>DRWR</b>	32.5	39.5	59.8	346.1
<b>DAWR</b>	38.6	46.0	70.5	354.3
<b>DIPO</b>	43.9	51.6	73.3	359.7
<b>IDQL</b>	33.8	41.7	63.7	349.9
<b>DQL</b>	36.9	44.4	68.5	353.5
<b>QSM</b>	64.4	72.4	107.6	391.2
<b>DPPo</b>	35.2	42.0	65.6	350.3

Table A2: **Wall-clock time** in seconds for a single training iteration in **ROBOMIMIC tasks with state input** when comparing diffusion-based RL algorithms. Each iteration involves 4 episodes (1200 environment timesteps for Lift and Can, 1600 for Square, and 3200 for Transport) from each of the 50 parallelized environments running on 50 CPU threads and a NVIDIA L40 GPU (60000, 80000, 160000 steps).

**Comparison of other policy parameterizations and architecture** **Section 5.2** and **Section 5.3.** Table A3 and Table A4 shows the wall-clock time used in fine-tuning in each ROBOMIMIC task with state or pixel input, respectively. Gaussian and GMM use similar times and Transformer is slightly more expensive than MLP. On average with state input, **DPPo**-MLP trains 24%, 21%, 24%, and 22% slower than baselines due to the more expensive diffusion sampling. **DPPo**-UNet requires more time with the extensive use of convolutional and normalization layers and trains on average 49% slower than **DPPo**-MLP. On average with pixel input, **DPPo**-ViT-MLP trains 14% slower than Gaussian-ViT-MLP — the difference is smaller than

the state input case as the rendering in simulation can be expensive. Table A5 shows the wall-clock time used in FURNITURE-BENCH tasks. **DPO**-UNet trains 20% slower than Gaussian-MLP on average.

Method	Task			
	Lift	Can	Square	Transport
Gaussian-MLP	27.7	35.7	56.2	255.6
Gaussian-Transformer	29.8	37.1	57.8	266.1
GMM-MLP	28.0	36.2	55.2	254.5
GMM-Transformer	29.5	37.4	58.1	260.2
<b>DPO</b> -MLP	35.6	43.3	65.0	350.5
<b>DPO</b> -UNet	83.6	92.7	130.4	431.1

Table A3: **Wall-clock time** in seconds for a single training iteration in **ROBOMIMIC tasks with state input** when comparing policy parameterizations. Each iteration involves 4 episodes (1200 environment timesteps for Lift and Can, 1600 for Square, and 3200 for Transport) from each of the 50 parallelized environments running on 50 CPU threads and a NVIDIA L40 GPU (60000, 80000, 160000 steps).

Method	Task			
	Lift	Can	Square	Transport
Gaussian-ViT-MLP	153.6	173.1	277.0	770.0
<b>DPO</b> -ViT-MLP	194.9	202.5	328.5	871.3

Table A4: **Wall-clock time** in seconds for a single training iteration in **ROBOMIMIC tasks with pixel input** when comparing policy parameterizations. Each iteration involves 4 episodes (1200 environment timesteps for Lift and Can, 1600 for Square, and 3200 for Transport) from each of the 50 parallelized environments running on 50 CPU threads and a NVIDIA L40 GPU (60000, 80000, 160000 steps).

Method	Task		
	One-leg	Lamp	Round-table
Gaussian-MLP	101.8	202.8	168.7
<b>DPO</b> -UNet	148.4	258.2	188.6

Table A5: **Wall-clock time** in seconds for a single training iteration in **FURNITURE-BENCH tasks** when comparing policy parameterizations. Each iteration involves 1 episodes (700 environment timesteps for One-leg, and 1000 for Lamp and Round-table) from each of the 1000 parallelized environments running on a NVIDIA L40 GPU (700000, 1000000, 1000000 steps).

## D Additional Experimental Details

### D.1 Details of policy architectures used in all experiments

**MLP.** For most of the experiments, we use a Multi-layer Perceptron (MLP) with two-layer residual connection as the policy head. For diffusion-based policies, we also use a small MLP encoder for the state input and another small MLP with sinusoidal positional encoding for the denoising timestep input. Their output features are then concatenated before being fed into the MLP head. Diffusion Policy, proposed by Chi et al.

[15], does not use MLP as the diffusion architecture, but we find it delivers comparable (or even better) pre-training performance compared to UNet.

**Transformer.** For comparing to other policy parameterizations in Section 5.2, we also consider Transformer as the policy architecture for the Gaussian and GMM baselines. We consider decoder only. No dropout is used. A learned positional embedding for the action chunk is the sequence into the decoder.

**UNet.** For comparing to other policy parameterizations in Section 5.2, we also consider UNet [66] as a possible architecture for DP. We follow the implementation from Chi et al. [15] that uses sinusoidal positional encoding for the denoising timestep input, except for using a larger MLP encoder for the observation input in each convolutional block. We find this modification helpful in more challenging tasks.

**ViT.** For pixel-based experiments in Section 5.2 we use Vision-Transformer(ViT)-based image encoder introduced by Hu et al. [31] before an MLP head. Proprioception input is appended to each channel of the image patches. We also follow [31] and use a learned spatial embedding for the ViT output to greatly reduce the number of features, which are then fed into the downstream MLP head.

Task	Obs dim - State	Obs dim - Pixel	Act dim	T	Sparse reward ?
GYM	Hopper-v2	11	-	3	1000
	Walker2D-v2	17	-	6	1000
	HalfCheetah-v2	17	-	6	1000
ROBOMIMIC, state input	Lift	19	-	7	300
	Can	23	-	7	300
	Square	23	-	7	400
	Transport	59	-	14	800
ROBOMIMIC, pixel input	Lift	9	96×96	7	300
	Can	9	96×96	7	300
	Square	9	96×96	7	400
	Transport	18	2×96×96	14	800
FURNITURE-BENCH	One-leg	58	-	10	700
	Lamp	44	-	10	1000
	Round-table	44	-	10	1000
D3IL	Avoid	4	-	2	100

Table A6: **Comparison of the different tasks considered.** “Obs dim - State”: dimension of the state observation input. “Obs dim - State”: dimension of the pixel observation input. “Act dim - State”: dimension of the action space.  $T$ : maximum number of steps in an episode. “Sparse reward ?”: whether sparse reward is used in training instead of dense reward.

## D.2 Additional details of GYM tasks and training in Section 5.1

**Pre-training.** The observations and actions are normalized to  $[0, 1]$  using min/max statistics from the pre-training dataset. For all three tasks the policy is trained for 3000 epochs with batch size 128, learning rate of 1e-3 decayed to 1e-4 with a cosine schedule, and weight decay of 1e-6. Exponential Moving Average (EMA) is applied with a decay rate of 0.995.

**Fine-tuning.** All methods from [Section 5.1](#) use the same pre-trained policy. Fine-tuning is done using online experiences sampled from 40 parallelized MuJoCo environments [[77](#)]. Reward curves shown in [Fig. 5](#) are evaluated by running fine-tuned policies with  $\sigma_{\min}^{\text{exp}} = 0.001$  (i.e., without extra noise) for 40 episodes. Each episode terminates if the default conditions are met or the episode reaches 1000 timesteps. Detailed hyperparameters are listed in [Table A7](#) and [Table A8](#).

### D.3 Descriptions of diffusion-based RL algorithm baselines in [Section 5.1](#)

**DRWR:** This is a **customized** reward-weighted regression (RWR) algorithm [[55](#)] that fine-tunes a pre-trained DP with a supervised objective with higher weights on actions that lead to higher reward-to-go  $r$ .

The reward is scaled with  $\beta$  and the exponentiated weight is clipped at  $w_{\max}$ . The policy is updated with experiences collected with the current policy (no buffer for data from previous iteration) and a replay ratio of  $N_\theta$ . No critic is learned.

$$\mathcal{L}_\theta = \mathbb{E}^{\bar{\pi}_\theta, \varepsilon_t} \left[ \min(e^{\beta r_t}, w_{\max}) \|\varepsilon_t - \varepsilon_\theta(a_t^0, s_t, k)\|^2 \right].$$

**DAWR:** This is a **customized** advantage-weighted regression (AWR) algorithm [[53](#)] that builds on **DRWR** but uses TD-bootstrapped [[75](#)] advantage estimation instead of the higher-variance reward-to-go for better training stability and efficiency. **DAWR** (and **DRWR**) can be seen as approximately optimizing ([4.2](#)) with a Kullback–Leibler (KL) divergence constraint on the policy [[53](#), [6](#)].

The advantage is scaled with  $\beta$  and the exponentiated weight is clipped at  $w_{\max}$ . Unlike **DRWR**, we follow [[53](#)] and trains the actor in an off-policy manner: recent experiences are saved in a replay buffer  $\mathcal{D}$ , and the actor is updated with a replay ratio of  $N_\theta$ .

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}, \varepsilon_t} \left[ \min(e^{\beta \hat{A}_\phi(s_t, a_t^0)}, w_{\max}) \|\varepsilon_t - \varepsilon_\theta(a_t^0, s_t, k)\|^2 \right].$$

The critic is updated less frequently (we find diffusion models need many gradient updates to fit the actions) with a replay ratio of  $N_\phi$ .

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|\hat{A}_\phi(s_t, a_t^0) - A(s_t, a_t^0)\|^2],$$

where  $A$  is calculated using  $\text{TD}(\lambda)$ , with  $\lambda$  as  $\lambda_{\text{DAWR}}$  and the discount factor  $\gamma_{\text{ENV}}$ .

**DIPO** [[86](#)]: This baseline applies “action gradient” that uses a learned state-action Q function to update the actions saved in the replay buffer, and then has DP fitting on them without weighting.

Similar to **DAWR**, recent experiences are saved in a replay buffer  $\mathcal{D}$ . The actions ( $k = 0$ ) in the buffer are updated for  $M_{\text{DIPO}}$  iterations with learning rate  $\alpha_{\text{DIPO}}$ .

$$a_t^{m+1, k=0} = a_t^{m, k=0} + \alpha_{\text{DIPO}} \nabla_\phi \hat{Q}_\phi(s_t, a_t^{m, k=0}), \quad m = 0, \dots, M_{\text{DIPO}} - 1.$$

The actor is then updated with a replay ratio of  $N_\theta$ .

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\|\varepsilon_t - \varepsilon_\theta(a_t^{M_{\text{DIPO}}, k=0}, s_t, k)\|^2].$$

The critic is trained to minimize the Bellman residual with a replay ratio of  $N_\phi$ . Double Q-learning is also applied.

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{Q}_\phi(s_{t+1}, \bar{\pi}_\theta(a_{t+1}^{k=0} | s_{t+1})) - \hat{Q}_\phi(s_t, a_t^{m=0, k=0})\|^2]$$

**IDQL [27]:** This baseline learns a state-action Q function and state V function to choose among the sampled actions from DP. DP fits on new samples without weighting.

Again recent experiences are saved in a replay buffer  $\mathcal{D}$ . The state value function is updated to match the expected Q value with an expectile loss, with a replay ratio of  $N_\psi$ .

$$\mathcal{L}_\psi = \mathbb{E}^{\mathcal{D}} [| \tau_{\text{IDQL}} - \mathbb{1}(\hat{Q}_\phi(s_t, a_t^0) < \hat{V}_\psi^2(s_t)) | ].$$

The value function is used to update the Q function with a replay ratio of  $N_\phi$ .

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\| (R_t + \gamma_{\text{ENV}} \hat{V}_\psi(s_{t+1}) - \hat{Q}_\phi(s_t, a_t^0))^2 \|].$$

The actor fits all sampled experiences without weighting, with a replay ratio of  $N_\theta$ .

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\| \varepsilon_t - \varepsilon_\theta(a_t^0, s_t, k) \|^2].$$

At inference time,  $M_{\text{IDQL}}$  actions are sampled from the actor. For training, Boltzmann exploration is applied based on the difference between  $Q$  value of the sampled actions and the  $V$  value at the current state. For evaluation, the greedy action under  $Q$  is chosen.

**DQL [84]:** This baseline learns a state-action Q function and backpropagates the gradient from the critic through the entire actor (with multiple denoising steps), akin to the usual Q-learning.

Again recent experiences are saved in a replay buffer  $\mathcal{D}$ . The actor is then updated using both a supervised loss and the value loss with a replay ratio of  $N_\theta$ .

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\| \varepsilon_t - \varepsilon_\theta(a_t^0, s_t, k) \|^2 - \alpha_{\text{DQL}} \hat{Q}_\phi(s_t, \bar{\pi}_\theta(a_t^0 | s_t))],$$

where  $\alpha_{\text{DQL}}$  is a weighting coefficient. The critic is trained to minimize the Bellman residual with a replay ratio of  $N_\phi$ . Double Q-learning is also applied.

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\| (R_t + \gamma_{\text{ENV}} \hat{Q}_\phi(s_{t+1}, \bar{\pi}_\theta(a_{t+1}^0 | s_{t+1})) - \hat{Q}_\phi(s_t, a_t^0))^2 \|]$$

**QSM [59]:** This baseline learns a state-action Q function, and then updates the actor by aligning the score of the diffusion actor with the gradient of the Q function.

Again recent experiences are saved in a replay buffer  $\mathcal{D}$ . The critic is trained to minimize the Bellman residual with a replay ratio of  $N_\phi$ . Double Q-learning is also applied.

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\| (R_t + \gamma_{\text{ENV}} \hat{Q}_\phi(s_{t+1}, \bar{\pi}_\theta(a_{t+1}^0 | s_{t+1})) - \hat{Q}_\phi(s_t, a_t^0))^2 \|].$$

The actor is updated as follows with a replay ratio of  $N_\theta$ .

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\| \alpha_{\text{QSM}} \nabla_a \hat{Q}_\phi(s_t, a_t) - (-\varepsilon_\theta(a_t^0, s_t, k)) \|^2],$$

where  $\alpha_{\text{QSM}}$  scales the gradient. The negative sign before  $\varepsilon_\theta$  is from taking the gradient of the mean  $\mu$  in the denoising process.

## D.4 Additional details of DPPO implementation in all tasks

Similar to all baselines in Appendix D.3, we denote  $N_\theta$  and  $N_\phi$  the replay ratio for the actor (Diffusion Policy) and the critic (state value function) in DPPO; in practice we always set  $N_\theta = N_\phi$  in DPPO, with the combined loss  $\mathcal{L} = \mathcal{L}_\theta + \mathcal{L}_\phi$ . Similar to usual PPO implementations [32], the batch updates in an iteration terminate when the KL divergence between  $\pi_\theta$  and  $\pi_{\theta_{\text{old}}}$  reaches 1.

We also find the PPO clipping ratio,  $\varepsilon$ , can affect the training stability significantly in DPPO (as well as in Gaussian and GMM policies) especially in sparse-reward manipulation tasks. In practice we find that, a good indicator of the amount of clipping leading to optimal training efficiency, is to aim for a clipping fraction (fraction of individual samples being clipped in a batch) of 10% to 20%. For each method in different tasks, we vary  $\varepsilon$  in  $\{.1, .01, .001\}$  and choose the highest value that satisfies the clipping fraction target. Empirically we also find that, using a higher  $\varepsilon$  for earlier denoising steps in DPPO further improves training stability in manipulation tasks. Denote  $\varepsilon_k$  the clipping value at denoising step  $k$ , and in practice we set  $\varepsilon_{k=(K-1)} = 0.1\varepsilon_{k=0}$ , and it follows an exponential schedule among intermediate  $k$ .

## D.5 Additional details of ROBOMIMIC tasks and training in Section 5.2

**Tasks.** We consider four tasks from the ROBOMIMIC benchmark [48]: (1) Lift: lifting a cube from the table, (2) Can: picking up a Coke can and placing it at a target bin, (3) Square: picking up a square nut and place it on a rod, and (4) Transport: two robot arms removing a bin cover, picking and placing a cube, and then transferring a hammer from one container to another one.

**Pre-training.** ROBOMIMIC provides the Multi-Human (MH) dataset with noisy human demonstrations for each task, which we use to pre-train the policies. The observations and actions are normalized to  $[0, 1]$  using min/max statistics from the pre-training dataset. No history observation (pixel, proprioception, or ground-truth object states) is used. All policies are trained with batch size 128, learning rate 1e-4 decayed to 1e-5 with a cosine schedule, and weight decay 1e-6. Diffusion-based policies are trained with 8000 epochs, while Gaussian and GMM policies are trained with 5000 epochs — we find diffusion models require more gradient updates to fit the data well.

**Fine-tuning.** Diffusion-based, Gaussian, and GMM pre-trained policies are then fine-tuned using online experiences sampled from 50 parallelized MuJoCo environments [77]. Success rate curves shown in Fig. 5, Fig. 6, and Fig. A8 are evaluated by running fine-tuned policies with  $\sigma_{\min}^{\text{exp}} = 0.001$  (i.e., without extra noise) for 50 episodes. Episodes terminates only when they reach maximum episode lengths (shown in Table A6). Detailed hyperparameters are listed in Table A9.

**Pixel training.** We use the wrist camera view in Lift and Can, the third-person camera view in Square, and the two robot shoulder camera views in Transport. Random-shift data augmentation is applied to the camera images during both pre-training and fine-tuning. Gradient accumulation is used in fine-tuning so that the same batch size (as in state-input training) can fit on the GPU. Detailed hyperparameters are listed in Table A10.

## D.6 Descriptions of policy parameterization baselines in Section 5.2

**Gaussian.** We consider unimodal Gaussian with diagonal covariance, the most commonly used policy parameterization in RL. The standard deviation for each action dimension,  $\sigma_{\text{Gau}}$ , is fixed during pre-training; we also tried to learn  $\sigma_{\text{Gau}}$  from the dataset but we find the training very unstable. During fine-tuning  $\sigma_{\text{Gau}}$  is learned starting from the same fixed value and also clipped between 0.01 and 0.2. Additionally we clip

the sampled action to be within 3 standard deviation from the mean. As discussed in [Appendix D.4](#), we choose the PPO clipping ratio  $\varepsilon$  based on the empirical clipping fraction in each task. This setup is also used in the FURNITURE-BENCH experiments. We note that we spend significant amount of efforts tuning the Gaussian baseline, and our results with it are some of the best known ones in RL training for long-horizon manipulation tasks (exceeding our initial expectations), e.g., reaching  $\sim 100\%$  success rate in Lamp with Low randomness.

**GMM.** We also consider Gaussian Mixture Model as the policy parameterization. We denote  $M_{\text{GMM}}$  the number of mixtures. The standard deviation for each action dimension in each mixture,  $\sigma_{\text{GMM}}$ , is also fixed during pre-training. Again during fine-tuning  $\sigma_{\text{GMM}}$  is learned starting from the same fixed value and also clipped between 0.01 and 0.2.

## D.7 Additional details of FURNITURE-BENCH tasks and training in [Section 5.3](#)

**Tasks.** We consider three tasks from the FURNITURE-BENCH benchmark [28]: (1) One-leg: assemble one leg of a table by placing the tabletop in the fixture corner, grasping and inserting the table leg, and screwing in the leg, (2) Lamp: place the lamp base in the fixture corner, grasp, insert, and screw in the light bulb, and finally place the lamp shade, (3) Round-table: place a round tabletop in the fixture corner, insert and screw in the table leg, and then insert and screw in the table base. See [Fig. A10](#) for the visualized rollouts in simulation.

**Pre-training.** The pre-training dataset is collected in the simulated environments using a SpaceMouse<sup>8</sup>, a 6 DoF input device. The simulator runs at 10Hz. At every timestep, we read off the state of the SpaceMouse as  $\delta \mathbf{a} = [\Delta x, \Delta y, \Delta z, \Delta \text{roll}, \Delta \text{pitch}, \Delta \text{yaw}]$ , which is converted to a quaternion before passed to the environment step and stored as the action alongside the current observation in the trajectory. If  $|\Delta \mathbf{a}_i| < \varepsilon \forall i$  for some small  $\varepsilon = 0.05$  defining the threshold for a no-op, we do not record any action nor pass it to the environment. Discarding no-ops is important for allowing the policies to learn from demonstrations effectively. When the desired number of demonstrations has been collected (typically 50), we process the actions to convert the delta actions stored from the SpaceMouse into absolute pose actions by applying the delta action to the current EE pose at each timestep.

The observations and actions are normalized to  $[-1, 1]$  using min/max statistics from the pre-training dataset. No history observation (proprioception or ground-truth object states) is used, i.e., only the current observation is passed to the policy. All policies are trained with batch size 256, learning rate 1e-4 decayed to 1e-5 with a cosine schedule, and weight decay 1e-6. Diffusion-based policies are trained with 8000 epochs, while Gaussian policies are trained with 3000 epochs. Gaussian policies can easily overfit the pre-trained dataset, while diffusion-based policies are more resilient. Gaussian policies also require a very large MLP ( $\sim 10$  million parameters) to fit the data well.

**Fine-tuning.** Diffusion-based and Gaussian pre-trained policies are then fine-tuned using online experiences sampled from 1000 parallelized IsaacGym environments [47]. Success rate curves shown in [Fig. 7](#) are evaluated by running fine-tuned policies with  $\sigma_{\min}^{\text{exp}} = 0.001$  (i.e., without extra noise) for 1000 episodes. Episodes terminate only when they reach maximum episode length (shown in [Table A6](#)). Detailed hyperparameters are listed in [Table A11](#). We find a smaller amount of exploration noise (we set  $\sigma_{\min}^{\text{exp}}$  and  $\sigma_{\text{Gau}}$  to be 0.04) is necessary for the pre-trained policy achieving nonzero success rates at the beginning of fine-tuning.

---

<sup>8</sup><https://3dconnexion.com/us/product/spacemouse-wireless/>

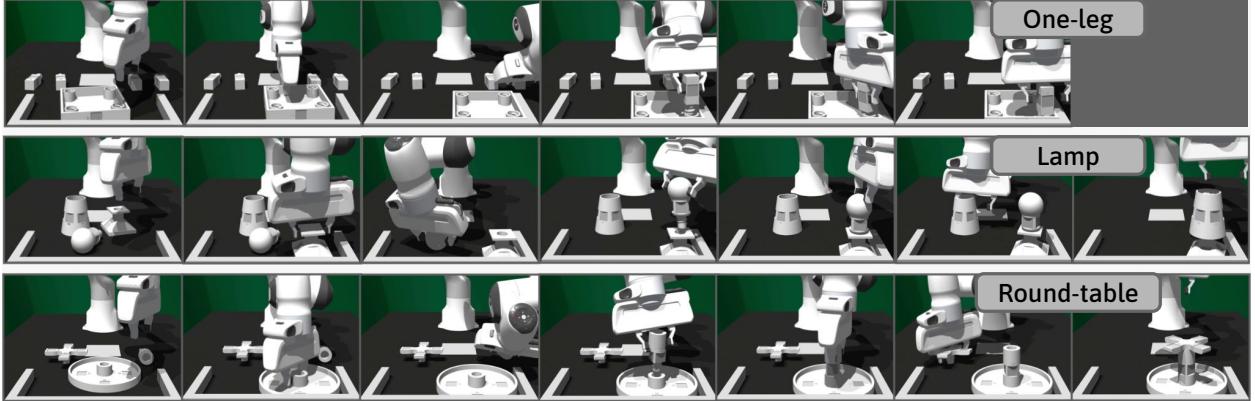


Figure A10: Representative rollouts from simulated FURNITURE-BENCH tasks.

**Hardware setup - robot control.** The physical robot used is a Franka Emika Panda arm. The policies output a sequence of desired end-effector poses in the robot base frame to control the robot. These poses are converted into joint position targets through differential inverse kinematics. We calculate the desired end-effector velocity as the difference between the desired and current poses divided by the delta time  $dt = 1/10$ . We then convert this to desired joint velocities using the Jacobian and compute the desired joint positions with a first-order integration over the current joint positions and desired velocity. The resulting joint position targets are passed to a low-level joint impedance controller provided by Polymetis [44], running at 1kHz.

**Hardware setup - state estimation.** To deploy state-based policies on real hardware, we utilize AprilTags [82] for part pose estimation. The FURNITURE-BENCH [28] task suite provides AprilTags for each part and code for estimating part poses from tag detections. The process involves several steps: (1) detecting tags in the camera frame, (2) mapping tag detections to the robot frame for policy compatibility, (3) utilizing known offsets between tags and object centers in the simulator, and (4) calibrating the camera pose using an AprilTag at a known position relative to the robot base. Despite general accuracy, detections can be noisy, especially during movement or partial occlusion. To mitigate occasional large “jumps” and noise in pose estimations, we implement a simple low-pass filter (exponential moving average) and average over multiple estimates from different tags or cameras for the same object.

**Hardware evaluation.** We perform 20 trials for each method. We adopt a single-blind model selection process: at the beginning of each trial, we first randomize the initial state. Then, we randomly select a method and roll it out, but the experimenter does not observe which model is used. We record the success and failure of each trial and then aggregate statistics for each model after all trials are completed.

**Domain randomization for sim-to-real transfer.** To facilitate the sim-to-real transfer, we apply additional domain randomization to the simulation training. We record the range of observation noises in hardware without any robot motion and then apply the same amount of noise to state observations in simulation. We find the state estimation in hardware particularly sensitive to the object heights. Also, we apply random noise (zero mean with 0.03 standard deviation) to the sampled action from **DPO** to simulate the imperfect low-level controller; we find adding such noise to the Gaussian policy leads to zero task success rate while **DPO** is robust to it (also see discussion in Section 6).

**BC regularization loss used for Gaussian baseline.** Since the fine-tuned Gaussian policy exhibits very jittery behavior and leads to zero success rate in real evaluation, we further experiment with adding a behav-

ior cloning (BC) regularization loss in fine-tuning with the Gaussian baseline. The combined loss follows

$$\mathcal{L}_{\theta,+\text{BC}} = \mathcal{L}_{\theta} - \alpha_{\text{BC}} \mathbb{E}^{\pi_{\theta_{\text{old}}}} \left[ \sum_{k=0}^{K-1} \log \pi_{\theta_{\text{pre-trained}}} (a_t^k | a_t^{k+1}, s_t) \right],$$

where  $\pi_{\theta_{\text{pre-trained}}}$  is the frozen BC-only policy. The extra term encourages the newly sampled actions from the fine-tuned policy to remain high-likelihood under the BC-only policy. We set  $\alpha_{\text{BC}} = 0.1$ . However, although this regularization reduces the sim-to-real gap, it also significantly limits fine-tuning, leading to the fine-tuning policy saturating at 53% success rate shown in Fig. 7.

## D.8 Additional details of **Avoid** task from D3IL and training in Section 6

**Pre-training.** We split the original dataset from D3IL based on the three settings, M1, M2, and M3; in each setting, observations and actions are normalized to  $[0, 1]$  using min/max statistics. All policies are trained with batch size 16 (due to the small dataset size), learning rate 1e-4 decayed to 1e-5 with a cosine schedule, and weight decay 1e-6. Diffusion-based policies are trained with about 15000 epochs, while Gaussian and GMM policies are trained with about 10000 epochs; we manually examine the trajectories from different pre-trained checkpoints and pick ones that visually match the expert data the best.

**Fine-tuning.** Diffusion-based, Gaussian, and GMM pre-trained policies are then fine-tuned using online experiences sampled from 50 parallelized MuJoCo environments [77]. Reward curves shown in Fig. 10 and Fig. A7 are evaluated by running fine-tuned policies with the same amount of exploration noise used in training for 50 episodes; we choose to use the training (instead of evaluation) setup since Gaussian policies exhibit multi-modality only with training noise. Episodes terminate only when they reach 100 steps.

**Added action noise during fine-tuning.** In Fig. 10 left, we demonstrate that **Dppo** exhibits stronger training stability when noise is added to the sampled actions during fine-tuning. The noise starts at the 5th iteration. It is sampled from a uniform distribution with the lower limit ramping up to 0.1 and the upper limit ramping up to 0.2 linearly in 5 iterations. The limits are kept the same from the 10th iteration to the end of fine-tuning.

## D.9 Listed training hyperparameters

Method	Parameter	Task(s)			
		GYM	Lift, Can	Square	Transport
Common	$\gamma_{\text{ENV}}$	0.99	0.999	0.999	0.999
	$\sigma_{\min}^{\text{exp}}$	0.1	0.1	0.1	0.08
	$\sigma_{\min}^{\text{prob}}$			0.1	
	$T_a$	4	4	4	8
	$K$			20	
	$K'$			20	
	Actor learning rate	1e-4	1e-5	1e-5	1e-5 (decayed to 1e-6)
	Critic learning rate (if applies)			1e-3	
<b>DRWR</b>	Actor MLP dims	[512, 512, 512]	[512, 512, 512]	[1024, 1024, 1024]	[1024, 1024, 1024]
	Critic MLP dims (if applies)			[256, 256, 256]	
	$\beta$			10	
<b>DAWR</b>	$w_{\max}$			100	
	$N_{\theta}$			16	
	Batch size			1000	
	$\beta$			10	
	$w_{\max}$			100	
<b>DIPO</b>	$\lambda_{\text{DAWR}}$			0.95	
	$N_{\theta}$			64	
	$N_{\phi}$	16	4	4	4
	Buffer size	200000	150000	150000	150000
	Batch size			256	
<b>IDQL</b>	$\alpha_{\text{DIPO}}$			1e-4	
	$M_{\text{DIPO}}$			10	
	$N_{\theta}$			64	
	Buffer size			400000	
	Batch size			5000	
<b>DQL</b>	$M_{\text{IDQL}}$	20	10	10	10
	$N_{\theta}$			16	
	$N_{\phi}$			16	
	Buffer size	200000	150000	150000	150000
	Batch size	256	512	512	512
	$\alpha_{\text{DQL}}$			1	
	$N_{\theta}$			64	
	$N_{\phi}$			64	
	Buffer size			400000	
	Batch size			5000	

Table A7: Fine-tuning hyperparameters for OpenAI GYM and ROBOMIMIC tasks when **comparing diffusion-based RL methods**. We list hyperparameters shared by all methods first, and then method-specific ones.

Method (cont'd)	Parameter (cont'd)	Task(s) (cont'd)			
		GYM	Lift, Can	Square	Transport
<b>QSM</b>	$\alpha_{\text{QSM}}$			50	
	$N_\theta$			32	
	$N_\phi$			32	
	Buffer size	200000	150000	150000	150000
<b>DPO</b>	Batch size			256	
	$\gamma_{\text{DENOISE}}$			0.99	
	GAE $\lambda$			0.95	
	$N_\theta$	5	10	10	8
	$N_\phi$	5	10	10	8
	$\varepsilon$			0.01	
	Batch size	5000	7500	10000	10000

Table A8: Continuation of Table A7.

Method	Parameter	Task		
		Lift, Can	Square	Transport
Common	$\gamma_{\text{ENV}}$		0.999	
	$T_a$	4	4	8
	Actor learning rate	1e-4	1e-5	1e-5 (decayed to 1e-6)
	Critic learning rate		1e-3	
	GAE $\lambda$		0.95	
	$N_\theta$	10	10	8
	$N_\phi$	10	10	8
	$\varepsilon$		0.01 (annealed in DPO)	
	Batch size	7500	10000	10000
Gaussian, Common	$\sigma_{\text{Gau}}$	0.1	0.1	0.08
Gaussian-MLP	Model size	552K	2.15M	1.93M
Gaussian-Transformer	Model size	675K	1.86M	1.87M
GMM, Common	$M_{\text{GMM}}$		5	
GMM-MLP	$\sigma_{\text{GMM}}$	0.1	0.1	0.08
	Model size	1.15M	4.40M	4.90M
	Model size	680K	1.87M	1.89M
<b>DPO</b> , Common	$\gamma_{\text{DENOISE}}$		0.99	
	$\sigma_{\min}^{\text{exp}}$	0.1	0.1	0.08
	$\sigma_{\min}^{\text{prob}}$	0.1	0.1	0.1
	$K$		20	
	$K'$		10	
DPO-MLP	Model size	576K	2.31M	2.43M
DPO-UNet	Model size	652K	1.62M	1.68M

Table A9: Fine-tuning hyperparameters for ROBOMIMIC tasks with **state** input when **comparing policy parameterizations**. We list hyperparameters shared by all methods first, and then method-specific ones. Since the different policy parameterizations use different neural network architecture, we list the total model size here instead of the details such as MLP dimensions.

Method	Parameter	Task		
		Lift, Can	Square	Transport
Common	$\gamma_{\text{ENV}}$		0.999	
	$T_a$	4	4	8
	Actor learning rate	1e-4	1e-5	1e-5 (decayed to 1e-6)
	Critic learning rate		1e-3	
	GAE $\lambda$		0.95	
	$N_\theta$	10	10	8
	$N_\phi$	10	10	8
	$\varepsilon$		0.01 (annealed in <b>DPO</b> )	
Gaussian-ViT-MLP	Batch size	7500	10000	10000
	Model size	1.03M	1.03M	1.93M
<b>DPO</b> -ViT-MLP	$\sigma_{\text{Gau}}$	0.1	0.1	0.08
	Model size	1.06M	1.06M	2.05M
	$\gamma_{\text{DENOISE}}$		0.9	
	$\sigma_{\text{min}}^{\text{exp}}$	0.1	0.1	0.08
	$\sigma_{\text{prob}}^{\text{min}}$		0.10	
	$K$		100	
	$K'$		5 (DDIM)	

Table A10: Fine-tuning hyperparameters for ROBOMIMIC tasks with **pixel** input when **comparing policy parameterizations**. We list hyperparameters shared by all methods first, and then method-specific ones. Since the different policy parameterizations use different neural network architecture, we list the total model size here instead of the details such as MLP dimensions.

Method	Parameter	Task		
		One-leg	Lamp	Round-table
Common	$\gamma_{\text{ENV}}$		0.999	
	$T_a$		8	
	Actor learning rate		1e-5 (decayed to 1e-6)	
	Critic learning rate		1e-3	
	GAE $\lambda$		0.95	
	$N_\theta$		5	
	$N_\phi$		5	
	$\varepsilon$		0.001	
Gaussian-MLP	Batch size		8800	
	Model size	10.64M	10.62M	10.62M
<b>DPO</b> -UNet	$\sigma_{\text{Gau}}$	0.04		
	Model size	6.86M	6.81M	6.81M
	$\gamma_{\text{DENOISE}}$		0.9	
	$\sigma_{\text{min}}^{\text{exp}}$		0.04	
	$\sigma_{\text{prob}}^{\text{min}}$		0.1	
	$K$		100	
	$K'$		5 (DDIM)	

Table A11: Fine-tuning hyperparameters for FURNITURE-BENCH tasks when **comparing policy parameterizations**. We list hyperparameters shared by all methods first, and then method-specific ones.