
DIFFUSION POLICY POLICY OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce *Diffusion Policy Policy Optimization*, **DPPO**, an algorithmic framework including best practices for fine-tuning diffusion-based policies (e.g. Diffusion Policy (Chi et al., 2024b)) 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 **DPPO** 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 **DPPO** 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 **DPPO** 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 videos: diffusionppoanon.github.io. Code: [anonymous.dppo](https://github.com/anonymous-dppo).

1 INTRODUCTION

Large-scale pre-training with additional fine-tuning has become a ubiquitous pipeline in the development of language and image foundation models (Brown et al., 2020; Radford et al., 2021; Ouyang et al., 2022; Ruiz et al., 2023). Though behavior cloning with expert data (Pomerleau, 1988) is rapidly emerging as dominant paradigm for pre-training *robot policies* (Florence et al., 2019; 2022; Zhao et al., 2023; Lee et al., 2024; Fu et al., 2024), their performance can be suboptimal (Osa et al., 2018) 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) (Sutton and Barto, 2018) 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 (Ho et al., 2020), which have emerged as a leading parameterization for action policies (Chi et al., 2024b; Reuss et al., 2023; Pearce et al., 2023), due in large part to their high training stability and ability to represent complex distributions (Rombach et al., 2022; Poole et al., 2022; Kong et al., 2020; Ho et al., 2022).

Contribution 1 (DPPO). We introduce *Diffusion Policy Policy Optimization (DPPO)*, 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 (Sutton et al., 1999; Schulman et al., 2017) in reinforcement learning.

The literature has already studied improving/fine-tuning diffusion-based policies (*Diffusion Policy*) using RL (Psenka et al., 2023; Wang et al., 2022; Hansen-Estruch et al., 2023), and has applied policy gradient (PG) to fine-tuning non-interactive applications of diffusion models such as text-to-image generation (Black et al., 2023; Clark et al., 2023; Fan et al., 2024). Yet PG methods have been believed to be inefficient in training Diffusion Policy for continuous control tasks (Psenka et al., 2023; Yang et al., 2023). 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.

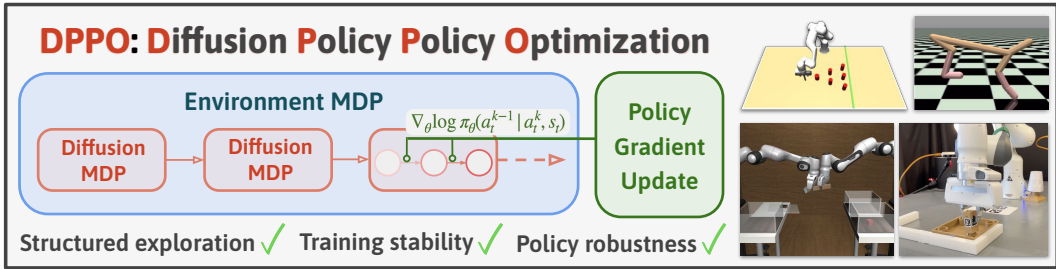


Figure 1: We introduce **DPPO**, *Diffusion Policy Policy Optimization*, that fine-tunes pre-trained Diffusion Policy using policy gradient. **DPPO** affords structured exploration and training stability during policy fine-tuning, and the fine-tuned policy exhibits strong robustness and generalization. **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 often final policy performance in comparison to a range of alternatives, including those based on off-policy Q-learning (Wang et al., 2022; Hansen-Estruch et al., 2023; Yang et al., 2023; Psenka et al., 2023) and weighted regression (Peng et al., 2019; Peters and Schaal, 2007; Kang et al., 2024), other demo-augmented RL methods (Ball et al., 2023; Nakamoto et al., 2024; Hu et al., 2023), 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 Janner et al. (2022); Ajay et al. (2023)). 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 show **DPPO** is effective in challenging robotic and control settings, including pixel observations and long-horizon manipulation tasks with sparse reward. We 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 (Black et al., 2023; Clark et al., 2023) to a multi-turn interactive setting with human feedback; drug design/discovery applications (Luo et al., 2022; Huang et al., 2024) with policy search on the molecular level in feedback with simulators (in the spirit of prior non-diffusion-based drug discovery with RL (Popova et al., 2018)); and the adaptation of diffusion-based language modeling (Sahoo et al., 2024; Lou et al., 2024) to interactive (e.g. with human feedback (Ouyang et al., 2022)), problem-solving and planning tasks.

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 (Williams, 1992; Sutton et al., 1999), there have been a range of algorithms that further improve training stability and sample efficiency such as DDPG (Lillicrap et al., 2015) and PPO (Schulman et al., 2015). PG methods have been broadly effective in training robot policies (Andrychowicz et al., 2020; Hwangbo

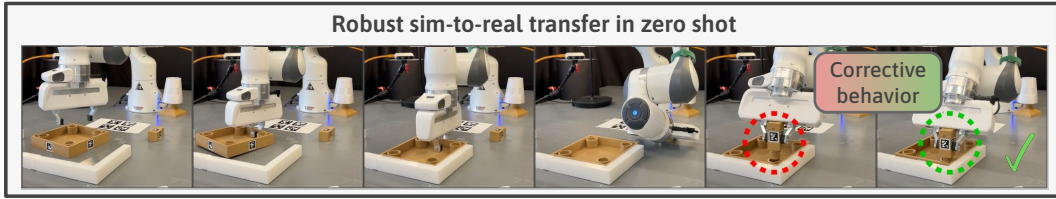


Figure 2: **DPPO** solves challenging long-horizon manipulation tasks from FURNITURE-BENCH (Heo et al. (2023)), enabling robust sim-to-real transfer without using any real data (Section 5.4).

et al., 2019; Kaufmann et al., 2023; Chen et al., 2022b), 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 (Rajeswaran et al., 2017; Torne et al., 2024; Peng et al., 2021). Our experiments find **DPPO** performing on-policy PG achieves stronger final performance in manipulation tasks, especially when the demonstrations are noisy, than off-policy Q-learning methods (Ball et al., 2023; Nakamoto et al., 2024; Hu et al., 2023). See Appendix A.1 for extended discussions on PG fine-tuning of robot policies and related methods.

Learning and improving diffusion-based policies. Diffusion-based policies (Chi et al., 2024b; Reuss et al., 2023; Ankile et al., 2024a; Ze et al.; Wang et al., 2024; Sridhar et al., 2023; Pearce et al., 2023) 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 (Janner et al., 2022; Ajay et al., 2023; Liang et al., 2023; Venkatraman et al., 2023; Chen et al., 2024). More recent work has explored techniques including Q-learning and weighted regression, either from purely offline estimation (Chen et al., 2022a; Wang et al., 2022; Ding and Jin, 2023), and/or with online interaction (Kang et al., 2024; Hansen-Estruch et al., 2023; Psenka et al., 2023; Yang et al., 2023). See Appendix A.2 for detailed descriptions of these methods.

Policy gradient through diffusion models. RL techniques have been used to fine-tune diffusion models such as ones for text-to-image generation (Fan and Lee, 2023; Fan et al., 2024; Black et al., 2023; Wallace et al., 2023). Black et al. (2023) treat the denoising process as an MDP and apply PPO updates. 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. (2023) 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 a *Markov Decision Process* (MDP)¹ $\mathcal{M}_{\text{ENV}} := (\mathcal{S}, \mathcal{A}, P_0, P, R)$ 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 $R(s_t, a_t)$ ². Fixing the MDP \mathcal{M}_{ENV} , we let \mathbb{E}^π (resp. \mathbb{P}^π) denote the expectation (resp. probability distribution) over trajectories $(s_0, a_0, \dots, s_T, a_T)$ 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 a cost $\mathcal{J}(\pi_\theta) = \mathbb{E}^{\pi_\theta, P_0}[\sum_{t \geq 0} \gamma(t)R(s_t, a_t)]$

¹More 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.

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

Policy optimization. The *policy gradient method* (e.g., REINFORCE (Williams, 1992)) allows for improving policy performance by approximating the gradient of this objective w.r.t. the policy parameters: $\nabla_{\theta} \mathcal{J}(\pi_{\theta}) = \mathbb{E}^{\pi_{\theta}, P_0} [\sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r_t(s_t, a_t)]$, $r_t(s_t, a_t) := \sum_{\tau \geq t} \gamma(\tau) R(s_{\tau}, a_{\tau})$ where r_t is the discounted cumulative future reward from time t (more generally, r_t can be replaced by a Q-function estimator (Sutton et al., 1999)), γ 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) (Nichol and Dhariwal, 2021; Ho et al., 2020; Sohl-Dickstein et al., 2015) 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 $\varepsilon_{\theta}(x_k, k)$, predicting the added noise ε that converts x_0 to x_k (Ho et al., 2020). Sampling starts with a random sample $x^K \sim \mathcal{N}(0, I)$ and iteratively generates the denoised sample:

$$x^{k-1} \sim p_{\theta}(x^{k-1} | x^k) := \mathcal{N}(x^{k-1}; \mu_k(x^k, \varepsilon_{\theta}(x^k, k)), \sigma_k^2 I). \quad (3.1)$$

Above, $\mu_k(\cdot)$ is a fixed function, independent of θ , that maps x^k and predicted ε_{θ} to the next mean, and σ_k^2 is a variance term that abides by a fixed schedule from $k = 1, \dots, K$. We refer the reader to Chan (2024) for an in-depth survey.

Diffusion models as policies. *Diffusion Policy* (DP; see Chi et al. (2024b)) is a policy π_{θ} parameterized by a DDPM which takes in s as a conditioning argument, and parameterizes $p_{\theta}(a^{k-1} | a^k, s)$ as in (3.1). DPs can be trained via behavior cloning by fitting the conditional noise prediction $\varepsilon_{\theta}(a^k, s, k)$ to predict the 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 diffusion and non-diffusion baselines use the same chunk size.

4 DPPO: DIFFUSION POLICY POLICY OPTIMIZATION

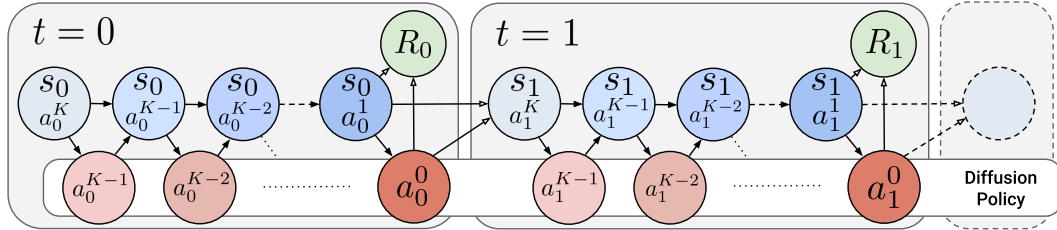


Figure 3: We treat the denoising process in Diffusion Policy as an MDP, and the whole environment episode can be considered as a chain of such MDPs. Now the entire chain (“Diffusion Policy MDP”, \mathcal{M}_{DP}) involves a Gaussian likelihood at each (denoising) step and thus can be optimized with policy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

As observed in Black et al. (2023) and Psenka et al. (2023), a denoising process can be represented as a multi-step MDP in which policy likelihood of each denoising step 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}_{DP} , visualized in Fig. 3. Below, we use the notation δ 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}_{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 $K - 1 \geq k \geq 0$. The states, actions and rewards are

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

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

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

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 (3.1), the policy in \mathcal{M}_{DP} is

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

Fortunately, (4.1) is a *Gaussian likelihood*, which can be evaluated analytically and is amenable to the policy gradient updates (see also Psenka et al. (2023) for an alternative derivation):

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

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.

4.1 INSTANTIATING **DPPO** WITH PROXIMAL POLICY OPTIMIZATION

We apply Proximal Policy Optimization (PPO) (Schulman et al., 2017; Engstrom et al., 2019; Huang et al., 2022; Achiam, 2018), a popular improvement of the vanilla policy gradient update.

Definition 4.1 (Generalized PPO, clipping variant). Consider a general MDP. Given an advantage estimator $\hat{A}(s, a)$, the PPO update (Schulman et al., 2017) is the sample approximation to

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

where ε , the clipping ratio, controls the maximum magnitude of the policy updated.

We instantiate PPO in our diffusion MDP with $(s, a, t) \leftarrow (\bar{s}, \bar{a}, \bar{t})$. Our advantage estimator 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 estimator³:

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

The denoising-discounting has the effect of downweighting the contribution of noisier steps (larger k) to the policy gradient (see study in Appendix D.2). Lastly, we choose the value estimator to *only depend* on the ‘‘s’’ component of \bar{s} : $\hat{V}^{\pi_{\theta_{\text{old}}}}(\bar{s}_{\bar{t}(t, 0)}) := \tilde{V}^{\pi_{\theta_{\text{old}}}}(s_t)$, 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 D.2.

Best Practices for **DPPO.** We summarize a number of best practices for **DPPO**; precise details are given in Appendix B. **(1)** We achieve substantial efficiency gains by fine-tuning the last few steps of the DDPM, whilst in many cases obtaining performance comparable to fine-tuning all steps. **(2)** For additional efficiency gains, one may fine-tune the Denoising Diffusion Implicit Model (DDIM) (Song et al., 2020a) instead. **(3)** We propose clipping the diffusion noise schedule at a larger-than-standard noise level to encourage exploration and training stability. **(4)** **DPPO** can be implemented with either Perceptron (MLP) or UNet (Ronneberger et al., 2015) as the policy head; while the former is simpler and achieves comparable performance to the latter, the latter provides flexibility in action chunk length which is useful in more challenging tasks.

³In practice, we use Generalized Advantage Estimation (GAE, Schulman et al. (2015)) that better balances variance and bias in estimating the advantage. We present the simpler form here.

5 PERFORMANCE EVALUATION OF DPPO

We study the performance of **DPPO** in popular RL and robotics benchmarking environments. Tables of all baselines and takeaways are given in Appendix C. While our evaluations focus primarily on **fine-tuning**, we also present training-from-scratch results in Appendix D. Wall-clock times are reported and discussed in Appendix E; they are roughly comparable (often faster) than other diffusion-based RL baselines, though can be up to $2\times$ slower than other policy parameterizations. Full choices of training hyperparameters and additional training details are presented in Appendix F.

Environments: OpenAI Gym. We first consider three population OpenAI GYM locomotion benchmarks (Brockman et al., 2016): {Hopper-v2, Walker2D-v2, HalfCheetah-v2}. All policies are pre-trained with the full medium-level datasets from D4RL (Fu et al., 2020) 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 (Mandlekar et al., 2021), {Lift, Can, Square, Transport}, ordered in increasing difficulty. These 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 Multi-Human demonstrations provided by ROBOMIMIC. We consider $T_a = 4$ for Can, Lift, and Square, and $T_a = 8$ for Transport. They are then fine-tuned with **sparse** reward upon task completion.

Environments: Furniture-Bench & real furniture assembly. Finally, we demonstrate solving longer-horizon, multi-stage robot manipulation tasks from the FURNITURE-BENCH (Heo et al., 2023) benchmark. We consider three simulated furniture assembly tasks, {One-leg, Lamp, Round-table}, shown in Fig. 4 and described in detail in Appendix F.8. We consider two levels of randomness over initial state distribution, Low and Med, defined by the benchmark. 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.

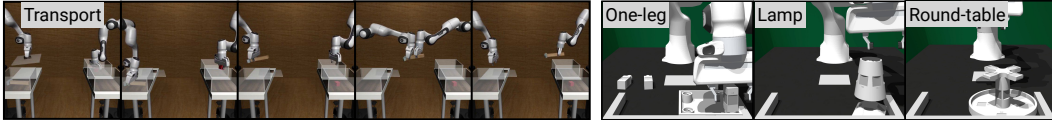


Figure 4: **Long-horizon robot manipulations tasks** including (left) the bimanual Transport from ROBOMIMIC and (right) FURNITURE-BENCH tasks (full rollouts visualized in Fig. 22).

5.1 COMPARISON TO DIFFUSION-BASED RL ALGORITHMS

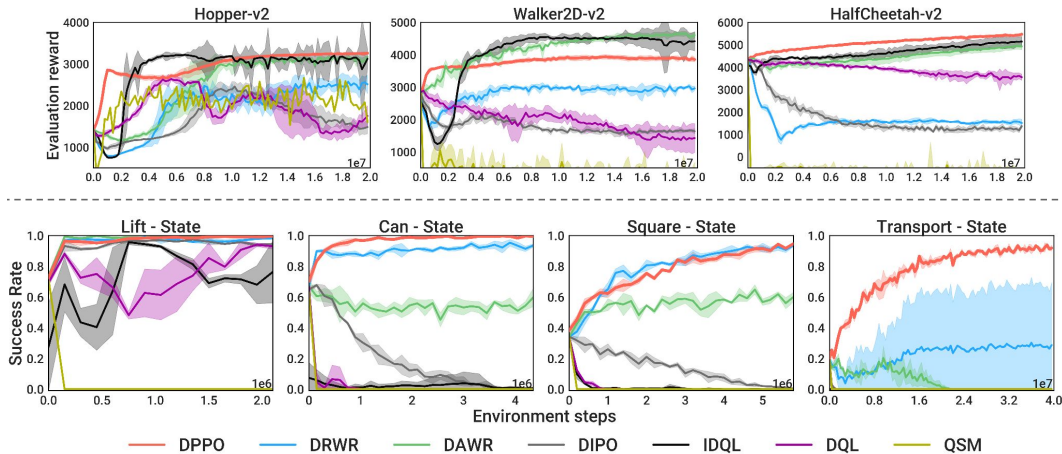
We compare **DPPO** to an extensive list of RL methods for fine-tuning diffusion models. **DRWR** and **DAWR** are *our own, novel* baselines based on reward-weighted regression (Peters and Schaal, 2007) and advantage-weighted regression (Peng et al., 2019). The remaining methods, **DIPO** (Yang et al., 2023), **IDQL** (Hansen-Estruch et al., 2023), **DQL** (Wang et al., 2022), and **QSM** (Psenka et al., 2023), are existing in the literature. We evaluate on the three OpenAI GYM tasks and the four ROBOMIMIC tasks with **state** input; see Appendix F.3 for further baseline and training details.

Overall, **DPPO** performs consistently, exhibits great training stability, and enjoys strong fine-tuning performance 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, using off-policy updates, suffers from training instability in sparse-reward ROBOMIMIC tasks given continuous action space plus large action chunk sizes (see further studies in Appendix D.2).

5.2 COMPARISON TO OTHER DEMO-AUGMENTED RL ALGORITHMS

We compare **DPPO** with recently proposed RL methods for training robot policies (not necessarily diffusion-based) leveraging offline data, including **RLPD** (Ball et al., 2023), **Cal-QL** (Nakamoto

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338



339
340
341
342
343
344
345
346
347
348
349

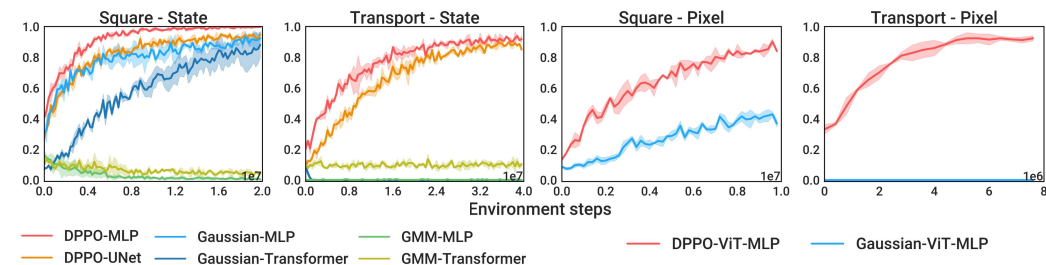
Figure 5: **Comparing to other diffusion-based RL algorithms.** Top row: GYM tasks (Brockman et al., 2016) averaged over five seeds. Bottom row: ROBOMIMIC tasks (Mandlekar et al., 2021), averaged over three seeds, with **state** observation. We opt to not show the erratic **QSM** error bar in Hopper-v2. **DPPO** curves are slightly thicker for better visibility.

et al., 2024), and **IBRL** (Hu et al., 2023). These methods add expert data in the replay buffer and performs off-policy updates (**IBRL** and **Cal-QL** also do pre-training), which significantly improves efficiency v.s. **DPPO** in HalfCheetah-v2. However, in sparse-reward manipulation tasks including Can and Square, **DPPO** achieves much better final performance than all three methods; **RLPD** and **Cal-QL** fail to learn at all and **IBRL** saturates at lower success levels. See Appendix D.1, containing Fig. 12, for further discussion.

350 5.3 COMPARISON TO OTHER POLICY PARAMETERIZATIONS

We compare **DPPO** with popular RL policy parameterizations: unimodal Gaussian with diagonal covariance (Sutton et al., 1999) and Gaussian Mixture Model (GMM, Bishop and Nasrabadi (2006)), using either MLPs or Transformers (Vaswani et al., 2017), 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.

360
361
362
363
364
365
366
367
368



369
370
371
372

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

373
374
375
376
377

Fig. 6 display results for the more challenging Square and Transport — we defer the results in Lift and Can to Fig. 20. With **state** input, **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. (2023) 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.**

5.4 EVALUATION ON FURNITURE-BENCH, AND SIM-TO-REAL TRANSFER

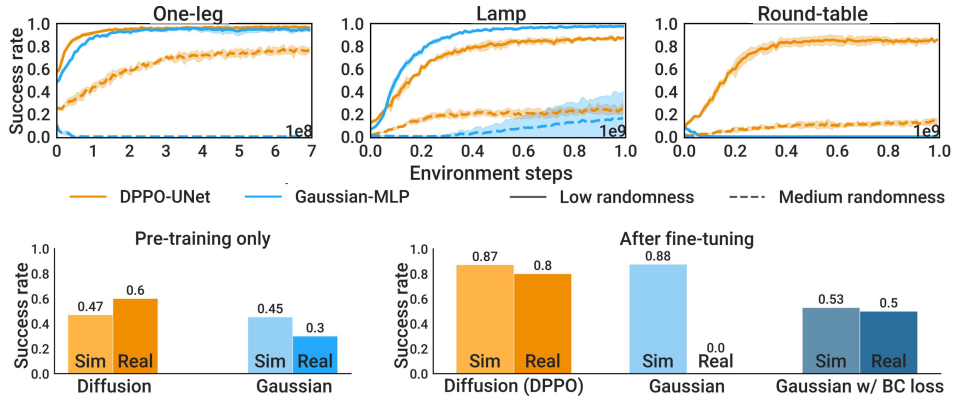


Figure 7: (Top) **DPPO** vs. Gaussian-MLP baseline in **simulated** FURNITURE-BENCH tasks. Results are averaged over three seeds. (Bottom) **Sim-to-real transfer results in One-leg**.

Here we evaluate **DPPO** on the long-horizon manipulation tasks from FURNITURE-BENCH (Heo et al., 2023). We compare **DPPO** to Gaussian-MLP, the overall most effective baseline from Section 5.3. Fig. 7 (top row) shows the evaluation success rate over fine-tuning iterations. **DPPO** 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 **DPPO** can leverage additional human data (better state space coverage) to further improve in Med, which is corroborated by ablation studies in Appendix D.3.

Sim-to-real transfer. We evaluate **DPPO** 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 F.8. Fig. 7 (bottom row) shows simulated and hardware success rates after pre-training and fine-tuning. Notably, **DPPO** 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 (Torne et al., 2024) 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; such behaviors are visualized in Fig. 2 and Fig. 23 shows representative rollouts on hardware.

5.5 SUMMARY OF ABLATION FINDINGS

Our ablation studies (c.f. Appendix D.2) find that: **(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 **DPPO** exploration, trading off between too little exploration and too much action noise; **(3)** **DPPO** is resilient to fine-tuning fewer-than- K denoising steps, yielding improved runtime and comparable performance; **(4)** **DPPO** 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 **DPPO**

We study the factors contributing to **DPPO**'s improvements in performance over the popular Gaussian and GMM methods introduced in Section 5.3. We use the `Avoid` environment from D3IL benchmark (Jia et al., 2024), where a robot arm needs to reach the other side of the table while avoiding an array of obstacles (Fig. 8, top-left). The action space is the 2D target location of the end-effector. D3IL provides expert demonstrations that covers different possible paths to the goal line — we consider three subsets of the demonstrations, M1, M2, and M3 in Fig. 8, each with two distinct modes; with only two modes in each setting, Gaussian (with exploration noise)⁴ 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 ($T_a = 4$ 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.

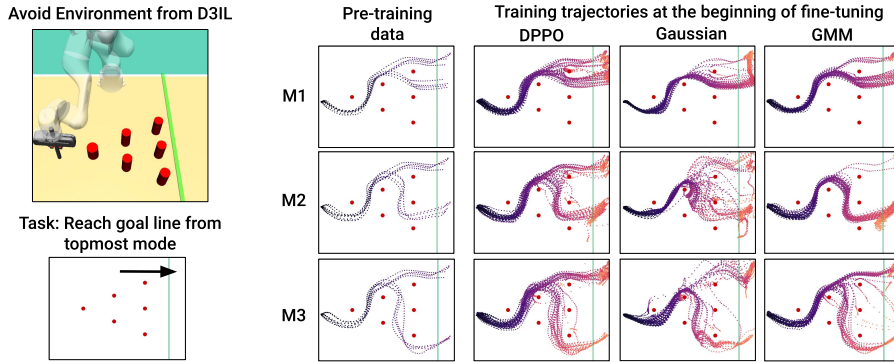


Figure 8: (Left) We use the `Avoid` environment from Jia et al. (2024) to visualize the **DPPO**'s exploration tendencies. The task is to reach the green goal line from the topmost mode. (Right) **Structured exploration.** We show sampled trajectories at the *first iteration of fine-tuning* for **DPPO**, Gaussian, and GMM after pre-training on three sets of expert demonstrations, M1, M2, and M3.

Benefit 1: Structured, on-manifold exploration. Fig. 8 (right) shows the sampled trajectories (with exploration noise) from **DPPO**, Gaussian, and GMM during the first iteration of fine-tuning. **DPPO** 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 **DPPO** is **structured in both action dimension and time horizon**.

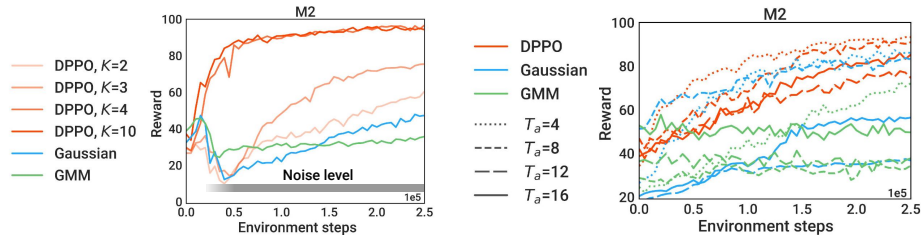
Benefit 2: Training stability from multi-step denoising process. In Fig. 9 (left), we run fine-tuning after pre-training with M2 and *attempt to de-stabilize fine-tuning* by gradually adding noise to the action during the fine-tuning process (see Appendix F.9 for details). We find that Gaussian and GMM's 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 5.4. In Fig. 9 (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.

Fig. 10 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**.

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. 11, we again add noise to the actions sampled from the fine-tuned policy (no noise applied during training) and find that **DPPO** policy exhibits strong robustness to the noise compared to the Gaussian policy. **DPPO** policy

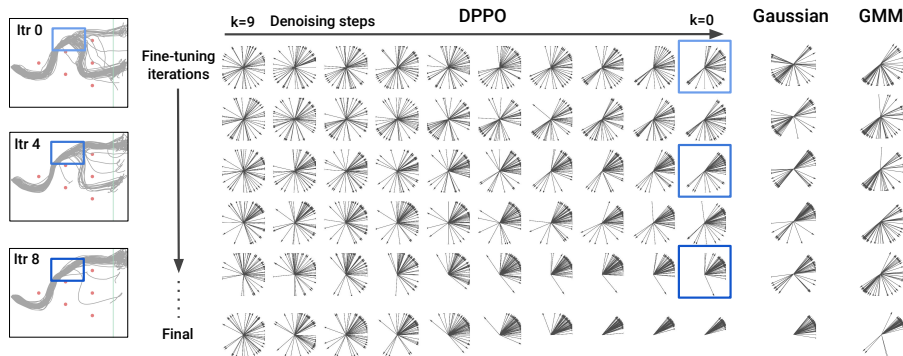
⁴Without noise, Gaussian policy is fully deterministic and cannot capture the two modes.

486
487
488
489
490
491
492
493



494 **Figure 9: Training stability.** Fine-tuning performance (averaged over five seeds, standard deviation
495 not shown) after pre-training with M2. (Left) Noise is injected into the applied actions after a few
496 training iterations. (Right) The action chunk size T_a is varied.

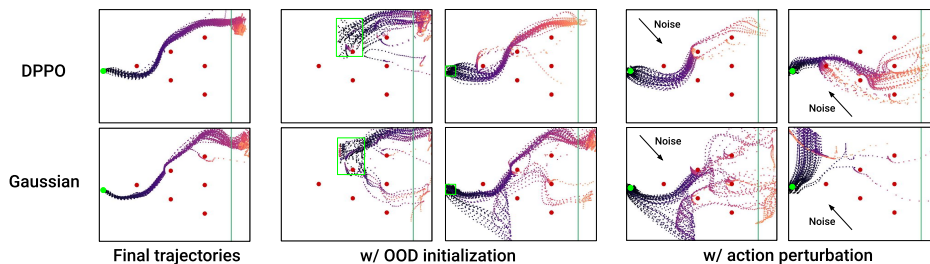
498
499
500
501
502
503
504
505
506
507
508
509



510 **Figure 10: Preserving the iterative refinement.** The 2D actions from 50 trajectories at the branching-
511 point *through fine-tuning* iterations after pre-training with M2. For **DPPO**, we also visualize the
512 action distribution through the final denoising steps at each fine-tuning iteration.

513 also converges to the (near-)optimal path from a larger distribution of initial states. This finding
514 echoes theoretical guarantees that Diffusion Policy, capable of representing complex multi-modal
515 data distribution, can effectively deconvolve noise from noisy states (Block et al., 2024), a property
516 used in Chen et al. (2024) to stabilize long-horizon video generation.

517
518
519
520
521
522
523
524
525



526 **Figure 11: Policy robustness after fine-tuning.** Green dot / box indicates the initial state region.

528

529 **7 CONCLUSION AND FUTURE WORK**

530
531
532
533
534
535
536
537
538
539

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 (Chen et al., 2023; Liang
et al., 2020; Ren et al., 2023; Chi et al., 2024a) 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 for 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 (Janner et al.,
2022; Ding et al., 2024) and decision-making aided by video prediction (Chen et al., 2024); these
tools may help address the main limitation of **DPPO** — its lower sample efficiency than off-policy
methods — and unlocking performing practical RL in physical hardware.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REFERENCES

- J. Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- 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.
- M. Alakuijala, G. Dulac-Arnold, J. Mairal, J. Ponce, and C. Schmid. Residual reinforcement learning from demonstrations. *arXiv preprint arXiv:2106.08050*, 2021.
- 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.
- L. Ankile, A. Simeonov, I. Shenfeld, and P. Agrawal. Juicer: Data-efficient imitation learning for robotic assembly. *arXiv*, 2024a.
- L. Ankile, A. Simeonov, I. Shenfeld, M. Torne, and P. Agrawal. From imitation to refinement-residual rl for precise visual assembly. *arXiv preprint arXiv:2407.16677*, 2024b.
- P. J. Ball, L. Smith, I. Kostrikov, and S. Levine. Efficient online reinforcement learning with offline data. In *International Conference on Machine Learning*, pages 1577–1594. PMLR, 2023.
- C. M. Bishop and N. M. Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.
- K. Black, M. Janner, Y. Du, I. Kostrikov, and S. Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- A. Block, A. Jadbabaie, D. Pfrommer, 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.
- G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 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.
- S. H. Chan. Tutorial on diffusion models for imaging and vision. *arXiv preprint arXiv:2403.18103*, 2024.
- 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.
- 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*, 2022a.
- T. Chen, J. Xu, and P. Agrawal. A system for general in-hand object re-orientation. In *Conference on Robot Learning*, 2022b.
- Y. Chen, C. Wang, L. Fei-Fei, and C. K. Liu. Sequential dexterity: Chaining dexterous policies for long-horizon manipulation. *arXiv preprint arXiv:2309.00987*, 2023.
- 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*, 2024a.
- 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*, 2024b.
- K. Clark, P. Vicol, K. Swersky, and D. J. Fleet. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.

594 Z. Ding and C. Jin. Consistency models as a rich and efficient policy class for reinforcement learning.
595 *arXiv preprint arXiv:2309.16984*, 2023.
596

597 Z. Ding, A. Zhang, Y. Tian, and Q. Zheng. Diffusion world model. *arXiv preprint arXiv:2402.03570*,
598 2024.
599

600 L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry. Implemen-
601 tation matters in deep rl: A case study on ppo and trpo. In *International conference on learning*
602 *representations*, 2019.

603 Y. Fan and K. Lee. Optimizing ddpm sampling with shortcut fine-tuning. *arXiv preprint*
604 *arXiv:2301.13362*, 2023.
605

606 Y. Fan, O. Watkins, Y. Du, H. Liu, M. Ryu, C. Boutilier, P. Abbeel, M. Ghavamzadeh, K. Lee,
607 and K. Lee. Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in*
608 *Neural Information Processing Systems*, 2024.

609 P. Florence, L. Manuelli, and R. Tedrake. Self-supervised correspondence in visuomotor policy
610 learning. *IEEE Robotics and Automation Letters*, 2019.
611

612 P. Florence, C. Lynch, A. Zeng, O. A. Ramirez, A. Wahid, L. Downs, A. Wong, J. Lee, I. Mordatch,
613 and J. Tompson. Implicit behavioral cloning. In *Conference on Robot Learning*. PMLR, 2022.

614 J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven rein-
615 forcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
616

617 Z. Fu, T. Z. Zhao, and C. Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost
618 whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024.
619

620 W. Goo and S. Niekum. Know your boundaries: The necessity of explicit behavioral cloning in
621 offline rl. *arXiv preprint arXiv:2206.00695*, 2022.

622 T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep
623 reinforcement learning with a stochastic actor. In *International conference on machine learning*,
624 pages 1861–1870. PMLR, 2018.
625

626 S. Haldar, J. Pari, A. Rai, and L. Pinto. Teach a robot to fish: Versatile imitation from one minute of
627 demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.

628 P. Hansen-Estruch, I. Kostrikov, M. Janner, J. G. Kuba, and S. Levine. Idql: Implicit q-learning as
629 an actor-critic method with diffusion policies. *arXiv preprint arXiv:2304.10573*, 2023.
630

631 M. Heo, Y. Lee, D. Lee, and J. J. Lim. Furniturebench: Reproducible real-world benchmark for
632 long-horizon complex manipulation. *arXiv preprint arXiv:2305.12821*, 2023.
633

634 T. Hester, M. Vecerik, O. Pietquin, M. Lanctot, T. Schaul, B. Piot, D. Horgan, J. Quan,
635 A. Sendonaris, I. Osband, et al. Deep q-learning from demonstrations. In *Proceedings of the*
636 *AAAI conference on artificial intelligence*, volume 32, 2018.

637 J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural infor-*
638 *mation processing systems*, 2020.
639

640 J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi,
641 D. J. Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv*
642 *preprint arXiv:2210.02303*, 2022.

643 H. Hu, S. Mirchandani, and D. Sadigh. Imitation bootstrapped reinforcement learning. *arXiv*
644 *preprint arXiv:2311.02198*, 2023.
645

646 S. Huang, R. F. J. Dossa, C. Ye, J. Braga, D. Chakraborty, K. Mehta, and J. G. AraÅšjo. Cleanrl:
647 High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of*
Machine Learning Research, 2022.

648 Z. Huang, L. Yang, X. Zhou, Z. Zhang, W. Zhang, X. Zheng, J. Chen, Y. Wang, C. Bin, and W. Yang.
649 Protein-ligand interaction prior for binding-aware 3d molecule diffusion models. In *The Twelfth*
650 *International Conference on Learning Representations*, 2024.

651 J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter. Learning
652 agile and dynamic motor skills for legged robots. *Science Robotics*, 2019.

653 M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior
654 synthesis. *arXiv preprint arXiv:2205.09991*, 2022.

655 X. Jia, D. Blessing, X. Jiang, M. Reuss, A. Donat, R. Lioutikov, and G. Neumann. Towards di-
656 verse behaviors: A benchmark for imitation learning with human demonstrations. *arXiv preprint*
657 *arXiv:2402.14606*, 2024.

658 B. Kang, X. Ma, C. Du, T. Pang, and S. Yan. Efficient diffusion policies for offline reinforcement
659 learning. *Advances in Neural Information Processing Systems*, 2024.

660 E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza. Champion-
661 level drone racing using deep reinforcement learning. *Nature*, 2023.

662 Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro. Diffwave: A versatile diffusion model for
663 audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.

664 S. Lee, Y. Wang, H. Etukuru, H. J. Kim, N. M. M. Shafiullah, and L. Pinto. Behavior generation
665 with latent actions. *arXiv preprint arXiv:2403.03181*, 2024.

666 J. Liang, S. Saxena, and O. Kroemer. Learning active task-oriented exploration policies for bridging
667 the sim-to-real gap. *arXiv preprint arXiv:2006.01952*, 2020.

668 Z. Liang, Y. Mu, M. Ding, F. Ni, M. Tomizuka, and P. Luo. Adaptdiffuser: Diffusion models as
669 adaptive self-evolving planners. *arXiv preprint arXiv:2302.01877*, 2023.

670 T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Contin-
671 uous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

672 Y. Lin, A. S. Wang, G. Sutanto, A. Rai, and F. Meier. Polymetis. [https://](https://facebookresearch.github.io/fairo/polymetis/)
673 facebookresearch.github.io/fairo/polymetis/, 2021.

674 A. Lou, C. Meng, and S. Ermon. Discrete diffusion modeling by estimating the ratios of the data
675 distribution. *stat*, 2024.

676 J. Luo, Z. Hu, C. Xu, Y. L. Tan, J. Berg, A. Sharma, S. Schaal, C. Finn, A. Gupta, and S. Levine.
677 Serl: A software suite for sample-efficient robotic reinforcement learning. *arXiv preprint*
678 *arXiv:2401.16013*, 2024.

679 S. Luo, Y. Su, X. Peng, S. Wang, J. Peng, and J. Ma. Antigen-specific antibody design and op-
680 timization with diffusion-based generative models for protein structures. *Advances in Neural*
681 *Information Processing Systems*, 2022.

682 V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,
683 A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for
684 robot learning. *arXiv preprint arXiv:2108.10470*, 2021.

685 A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu,
686 and R. Martín-Martín. What matters in learning from offline human demonstrations for robot
687 manipulation. In *arXiv preprint arXiv:2108.03298*, 2021.

688 A. Nair, A. Gupta, M. Dalal, and S. Levine. Awac: Accelerating online reinforcement learning with
689 offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.

690 M. Nakamoto, S. Zhai, A. Singh, M. Sobol Mark, Y. Ma, C. Finn, A. Kumar, and S. Levine. Cal-ql:
691 Calibrated offline rl pre-training for efficient online fine-tuning. *Advances in Neural Information*
692 *Processing Systems*, 36, 2024.

-
- 702 A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *International*
703 *conference on machine learning*, 2021.
- 704
- 705 T. Osa, J. Pajarinen, G. Neumann, J. A. Bagnell, P. Abbeel, J. Peters, et al. An algorithmic perspec-
706 tive on imitation learning. *Foundations and Trends® in Robotics*, 2018.
- 707
- 708 L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
709 K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback.
710 *Advances in neural information processing systems*, 2022.
- 711
- 712 T. Pearce, T. Rashid, A. Kanervisto, D. Bignell, M. Sun, R. Georgescu, S. V. Macua, S. Z. Tan,
713 I. Momennejad, K. Hofmann, et al. Imitating human behaviour with diffusion models. *arXiv*
714 *preprint arXiv:2301.10677*, 2023.
- 715
- 716 X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and
717 scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 718
- 719 X. B. Peng, Z. Ma, P. Abbeel, S. Levine, and A. Kanazawa. Amp: Adversarial motion priors for
720 stylized physics-based character control. *ACM Transactions on Graphics (ToG)*, 2021.
- 721
- 722 J. Peters and S. Schaal. Reinforcement learning by reward-weighted regression for operational space
723 control. In *Proceedings of the 24th international conference on Machine learning*, 2007.
- 724
- 725 D. A. Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in neural*
726 *information processing systems*, 1988.
- 727
- 728 B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion.
729 *arXiv preprint arXiv:2209.14988*, 2022.
- 730
- 731 M. Popova, O. Isayev, and A. Tropsha. Deep reinforcement learning for de novo drug design.
732 *Science advances*, 2018.
- 733
- 734 M. Psenka, A. Escontrela, P. Abbeel, and Y. Ma. Learning a diffusion model policy from rewards
735 via q-score matching. *arXiv preprint arXiv:2312.11752*, 2023.
- 736
- 737 A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
738 J. Clark, et al. Learning transferable visual models from natural language supervision. In *Inter-*
739 *national conference on machine learning*, 2021.
- 740
- 741 A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learning
742 complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv*
743 *preprint arXiv:1709.10087*, 2017.
- 744
- 745 A. Z. Ren, H. Dai, B. Burchfiel, and A. Majumdar. Adaptsim: Task-driven simulation adaptation
746 for sim-to-real transfer. In *Proceedings of the Conference on Robot Learning (CoRL)*, 2023.
- 747
- 748 M. Reuss, M. Li, X. Jia, and R. Lioutikov. Goal-conditioned imitation learning using score-based
749 diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 750
- 751 M. Rigter, J. Yamada, and I. Posner. World models via policy-guided trajectory diffusion. *arXiv*
752 *preprint arXiv:2312.08533*, 2023.
- 753
- 754 R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
755 with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and*
pattern recognition, 2022.
- 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.
- 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.
- 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.

756 J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control
757 using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
758

759 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
760 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

761 J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning
762 using nonequilibrium thermodynamics. In *International conference on machine learning*, 2015.
763

764 J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint*
765 *arXiv:2010.02502*, 2020a.

766 Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based genera-
767 tive modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020b.
768

769 A. Sridhar, D. Shah, C. Glossop, and S. Levine. Nomad: Goal masked diffusion policies for navi-
770 gation and exploration. *arXiv preprint arXiv:2310.07896*, 2023.
771

772 R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

773 R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement
774 learning with function approximation. *Advances in neural information processing systems*, 1999.
775

776 E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*
777 *IEEE/RSJ international conference on intelligent robots and systems*, 2012.

778 M. Torne, A. Simeonov, Z. Li, A. Chan, T. Chen, A. Gupta, and P. Agrawal. Reconciling real-
779 ity through simulation: A real-to-sim-to-real approach for robust manipulation. *arXiv preprint*
780 *arXiv:2403.03949*, 2024.
781

782 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Pol-
783 osukhin. Attention is all you need. *Advances in neural information processing systems*, 2017.
784

785 M. Vecerik, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and
786 M. Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems
787 with sparse rewards. *arXiv preprint arXiv:1707.08817*, 2017.

788 S. Venkatraman, S. Khaitan, R. T. Akella, J. Dolan, J. Schneider, and G. Berseth. Reasoning with
789 latent diffusion in offline reinforcement learning. *arXiv preprint arXiv:2309.06599*, 2023.
790

791 B. Wallace, M. Dang, R. Rafailov, L. Zhou, A. Lou, S. Purushwalkam, S. Ermon, C. Xiong, S. Joty,
792 and N. Naik. Diffusion model alignment using direct preference optimization. *arXiv preprint*
793 *arXiv:2311.12908*, 2023.

794 J. Wang and E. Olson. Apriltag 2: Efficient and robust fiducial detection. In *IEEE/RSJ International*
795 *Conference on Intelligent Robots and Systems (IROS)*, 2016.
796

797 L. Wang, J. Zhao, Y. Du, E. H. Adelson, and R. Tedrake. Poco: Policy composition from and for
798 heterogeneous robot learning. *arXiv preprint arXiv:2402.02511*, 2024.

799 Z. Wang, J. J. Hunt, and M. Zhou. Diffusion policies as an expressive policy class for offline
800 reinforcement learning. *arXiv preprint arXiv:2208.06193*, 2022.
801

802 R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
803 learning. *Machine learning*, 1992.

804 L. Yang, Z. Huang, F. Lei, Y. Zhong, Y. Yang, C. Fang, S. Wen, B. Zhou, and Z. Lin. Pol-
805 icy representation via diffusion probability model for reinforcement learning. *arXiv preprint*
806 *arXiv:2305.13122*, 2023.
807

808 Y. Ze, G. Zhang, K. Zhang, C. Hu, M. Wang, and H. Xu. 3d diffusion policy: Generalizable
809 visuomotor policy learning via simple 3d representations. In *ICRA 2024 Workshop on 3D Visual*
Representations for Robot Manipulation.

810 T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation with
811 low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
812

813 H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar. Dexterous manipulation with deep
814 reinforcement learning: Efficient, general, and low-cost. In *2019 International Conference on*
815 *Robotics and Automation (ICRA)*, pages 3651–3657. IEEE, 2019.

816 Z. Zhu, H. Zhao, H. He, Y. Zhong, S. Zhang, Y. Yu, and W. Zhang. Diffusion models for reinforce-
817 ment learning: A survey. *arXiv preprint arXiv:2311.01223*, 2023.
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864	CONTENTS	
865		
866	1 Introduction	1
867		
868	2 Related Work	2
869		
870	3 Preliminaries	3
871		
872		
873	4 DPPO: Diffusion Policy Policy Optimization	4
874	4.1 Instantiating DPPO with Proximal Policy Optimization	5
875		
876		
877	5 Performance Evaluation of DPPO	6
878	5.1 Comparison to diffusion-based RL algorithms	6
879	5.2 Comparison to other demo-augmented RL algorithms	6
880	5.3 Comparison to other policy parameterizations	7
881	5.4 Evaluation on Furniture-Bench, and sim-to-real transfer	8
882	5.5 Summary of ablation findings	8
883		
884		
885		
886	6 Understanding the performance of DPPO	9
887		
888		
889	7 Conclusion and Future Work	10
890		
891	A Extended Related Work	18
892	A.1 RL training of robot policies with offline data	18
893	A.2 Diffusion-based RL methods	18
894		
895		
896	B Best Practices for DPPO	19
897		
898	C Summary of All Baselines	20
899		
900	D Additional experimental results	21
901	D.1 Comparing to other demo-augmented RL methods	21
902	D.2 Ablation studies on design decisions in DPPO	21
903	D.3 Effect of expert data	23
904	D.4 Comparing to other policy parameterizations in <code>Avoid</code>	24
905	D.5 Comparing to other policy parameterizations in the easier tasks from <code>ROBOMIMIC</code>	24
906	D.6 Comparing to policy gradient using exact likelihood of Diffusion Policy	25
907		
908		
909		
910		
911	E Reporting of Wall-Clock Times	25
912		
913	F Additional Experimental Details	27
914	F.1 Details of policy architectures used in all experiments	27
915	F.2 Additional details of GYM tasks and training in Section 5.1	27
916	F.3 Descriptions of diffusion-based RL algorithm baselines in Section 5.1	28
917		

918	F.4	Descriptions of RL fine-tuning algorithm baselines in Section 5.2	29
919	F.5	Additional details of DPPO implementation in all tasks	30
920	F.6	Additional details of ROBOMIMIC tasks and training in Section 5.3	31
921	F.7	Descriptions of policy parameterization baselines in Section 5.3	31
922	F.8	Additional details of FURNITURE-BENCH tasks and training in Section 5.4	31
923	F.9	Additional details of <code>AVOID</code> task from D3IL and training in Section 6	35
924	F.10	Listed training hyperparameters	36

A EXTENDED RELATED WORK

A.1 RL TRAINING OF ROBOT POLICIES WITH OFFLINE DATA

Here, we discuss related work in training robot policies using RL augmented with offline data to help RL better explore online in sparse reward settings.

One simple form is to use offline data to pre-train the policy, typically using behavior cloning, and then fine-tune the policy online. This is the approach that **DPPO** takes. Often, a regularization loss is applied to constrain the fine-tuned policy to stay close to the base policy, leading to natural fine-tuned behavior and often better learning (Rajeswaran et al., 2017; Zhu et al., 2019; Torne et al., 2024). **DPPO** does not apply regularization at fine-tuning as we find the on-manifold exploration helps **DPPO** maintain natural behavior after fine-tuning Section 5.4. Another popular approach is to learn a *residual* policy with RL on top of the frozen base policy (Alakuijala et al., 2021; Haldar et al., 2023). A closer work to ours is Ankile et al. (2024b), which trains a one-step residual non-diffusion policy with on-policy RL on top of a pre-trained chunked diffusion policy. This approach has the benefit of being fully closed-loop but lacks the structured on-manifold exploration of **DPPO**. Another hybrid approach is from Hu et al. (2023), which uses pre-trained and fine-tuned policies to sample online experiences.

Another popular line of work, instead of training a base policy using offline data, directly adds the data in the replay buffer for online, off-policy learning in a single stage (Hester et al., 2018; Vecerik et al., 2017; Nair et al., 2020). One recent approach from Ball et al. (2023), **RLPD**, further improves sample efficiency from previous off-policy methods incorporating, e.g., critic ensembling. Luo et al. (2024) demonstrates **RLPD** solving real-world manipulation tasks (although generally less challenging than ones solved by **DPPO**). Other approaches including **Cal-Ql** build on offline RL to learn from offline data and then switch to online RL (Nakamoto et al., 2024; Hansen-Estruch et al., 2023). **IBRL** from Hu et al. (2023) pre-trains the base policy and samples offline data in fine-tuning.

A.2 DIFFUSION-BASED RL METHODS

This section discusses related methods that directly train or improve diffusion-based policies with RL 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. (2023) for an extensive survey on diffusion models for RL.

Most previous works 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. (2022) applies classifier guidance that generates trajectories with higher predicted rewards. Ajay et al. (2023) introduces classifier-free guidance that avoids learning the value of noisy states. There is another line of work that uses diffusion models as an action policy (instead of state planner) and generally applies Q-learning. **DQL** (Wang et al., 2022) 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) denoising chain, akin to the usual Q-learning. **IDQL** (Hansen-Estruch et al., 2023), 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. (2024) 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 (2022) similarly extracts the policy in the spirit of AWR (Peng et al., 2019). Chen et al. (2022a) 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. A small amount of work also focuses entirely on the online setting. **DIPO** (Yang et al., 2023) differs from **DQL** and related work in 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 (Psenka et al., 2023), 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. Rigger et al. (2023) proposes learning a diffusion dynamic model to generate synthetic trajectories for online training of a non-diffusion RL policy.

B BEST PRACTICES FOR **DPPO**

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 **DPPO**, we can choose to fine-tune only a subset of the denoising steps instead, e.g., the last K' steps. In Section 4.1 and Section 6 we find this speeds up **DPPO** training and reduces GPU memory usage without sacrificing the asymptotic performance. Instead of fine-tuning the pre-trained model weights θ , we make a copy θ_{FT} — θ is frozen and used for the early denoising steps, while θ_{FT} is used for the last K' steps and updated with **DPPO**.

Fine-tune DDIM. Instead of fine-tuning all K or the last few steps of the DDPM, one can also apply Denoising Diffusion Implicit Model (DDIM) (Song et al., 2020a) 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 **DPPO** efficiency as fewer steps are fine-tuned.

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

Although DDIM is typically used as a deterministic sampler by setting $\eta = 0$ in (B.1), we can use $\eta > 0$ for fine-tuning that provides exploration noise and avoids calculating Gaussian likelihood with a Dirac distribution. In practice, we set $\eta = 1$ for training (equivalent to applying DDPM Song et al. (2020a)) and then $\eta = 0$ for evaluation. *We reserve 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 σ_k introduced in Nichol and Dhariwal (2021), which was originally annealed to a small value on the order of $1E-4$ at $k = 0$. In **DPPO**, the value of σ_k also translates to the exploration noise that is crucial to training efficiency. Empirically, we find that clipping σ_k to a higher minimum value (denoted $\sigma_{\text{min}}^{\text{exp}}$, e.g., 0.01 – 0.1) when sampling actions helps exploration (see sensitivity analysis in Appendix D.2). Additionally we clip σ_k to be at least 0.1 (denoted $\sigma_{\text{min}}^{\text{prob}}$) 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.

Network architecture. We study both Multi-layer Perceptron (MLP) and UNet (Ronneberger et al., 2015) as the policy heads in Diffusion Policy. An MLP offers simpler setup and we find it generally fine-tunes more stably with **DPPO**. Moreover, since the UNet applies convolution to the denoised action, we can pre-train and fine-tune with different action chunk size T_a (the number of environment timesteps that the policy predicts future actions with), 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)⁵.

⁵With 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.

C SUMMARY OF ALL BASELINES

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Comparison to Other Diffusion RL Methods

<i>Method Name</i>	<i>Summary</i>
DPPO (ours)	Competitive on GYM; much stronger on ROBOMIMIC; the only one to solve Transport.
DAWR (ours)	Competitive on GYM; much weaker on ROBOMIMIC.
DRWR (ours)	Weaker on GYM; competitive on all of ROBOMIMIC but Transport
IDQL (Hansen-Estruch et al., 2023)	Competitive on GYM, much weaker on ROBOMIMIC.
DQL (Wang et al., 2022)	Much weaker on GYM and ROBOMIMIC.
QSM (Psenka et al., 2023)	Much weaker on GYM and ROBOMIMIC.
DIPO (Yang et al., 2023)	Much weaker on GYM and ROBOMIMIC.

Comparison to Other Demonstration-Augmented RL Methods

<i>Method Name</i>	<i>Summary</i>
DPPO (ours)	Much stronger on ROBOMIMIC; underperforms RLPD and Cal-QL on HalfCheetah-v2.
RLPD (Ball et al., 2023)	Very efficient on HalfCheetah-v2; zero reward on ROBOMIMIC.
Cal-QL (Nakamoto et al., 2024)	More efficient than DPPO but less efficient than RLPD on HalfCheetah-v2; zero reward on ROBOMIMIC.
IBRL (Hu et al., 2023)	Weaker than DPPO on ROBOMIMIC

Comparison to Other Policy Parameterization/Architecture

ROBOMIMIC State	<i>Summary</i>
DPPO-MLP (ours)	Performs best overall ; attains max reward on Square.
DPPO-UNet (ours)	Slightly underperforms DPPO-MLP ; second best.
Gaussian-MLP	Competitive on Square; zero reward on Transport.
Gaussian-Transformer	Weaker on Square; zero reward on Transport.
GMM-MLP	Low reward on Square; zero reward on Transport.
GMM-Transformer	Low reward on Square and Transport.
ROBOMIMIC Pixel	
DPPO-ViT-MLP (ours)	Performs better overall ; attains strong reward on Square and Transport.
Gaussian-ViT-MLP	Much weaker on Square; zero reward on Transport.
FURNITURE-BENCH	
DPPO-UNet (ours)	Performs better overall except slightly weaker on Lamp (Low randomness) and tied on One-leg Low.
Gaussian-MLP	Slightly stronger on Lamp with Low randomness and tied on One-leg Low; much weaker otherwise.
Sim-to-Real	
DPPO-UNet (ours)	Tied with Gaussian-MLP in sim; much stronger in transfer to real.
Gaussian-MLP	Strong in sim; zero success in real
Gaussian w/ BC Loss	Markedly weaker in sim; markedly weaker than DPPO (but non-zero reward) in real.

D ADDITIONAL EXPERIMENTAL RESULTS

D.1 COMPARING TO OTHER DEMO-AUGMENTED RL METHODS

In Section 5.2 we discuss that **DPPO** leads to superior final performance in manipulation tasks compared to other RL methods leveraging offline data, including **RLPD** (Ball et al., 2023), **Cal-QL** (Nakamoto et al., 2024), and **IBRL** (Hu et al., 2023). The full results are shown in Fig. 12 below. We use action chunk size $T_a = 1$ following the setup from these methods (**DPPO** may benefit from longer chunk size, albeit). The three baselines all achieve high reward in `HalfCheetah-v2` with much higher sample efficiency thanks to performing off-policy updates. However, in sparse-reward ROBOMIMIC tasks including `Can` and `Square`, **DPPO** outperforms all three significantly and achieves $\sim 100\%$ final success rates. **RLPD** and **Cal-QL** fail to achieve any success (0%) during evaluation, while **IBRL** saturates at lower success levels.

Our results with **RLPD** in `Can` and `Square` corroborates those from Hu et al. (2023). **IBRL** is shown to achieve high success rates ($> 90\%$) in both tasks in Hu et al. (2023); we hypothesize here it underperforms possibly due to (1) the noisier expert data (Multi-Human dataset from ROBOMIMIC) affects gradient update with mixed batches of online and offline data, and (2) our setup not using any history observation unlike Hu et al. (2023) stacking three observations.

Lastly, we note that although **DPPO** uses more environment steps, it runs significantly faster than the baselines as it leverages sampling from highly parallelized environments (40 in `HalfCheetah-v2` and 50 in `Can` and `Square`), while off-policy methods may fail to fully leverage such parallelized setup as the policy is updated less often and the performance may be affected.

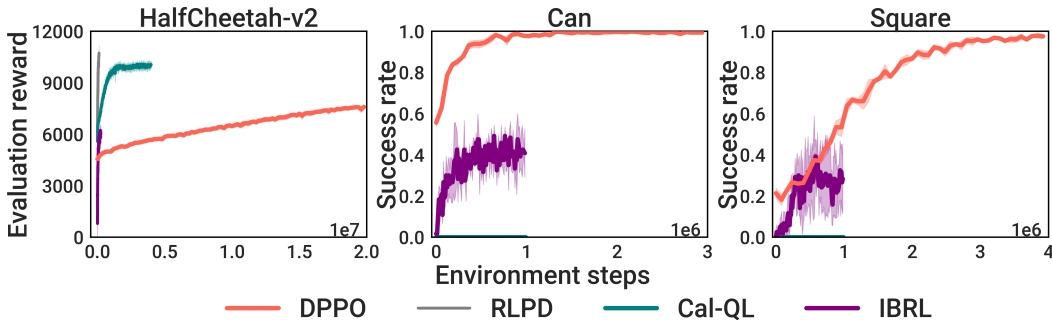


Figure 12: **Comparing to other demo-augmented RL methods.** Results are averaged over five seeds in `HalfCheetah-v2` and three seeds in `Can` and `Square`.

D.2 ABLATION STUDIES ON DESIGN DECISIONS IN **DPPO**

1. Choice of advantage estimator. In Section 4.1 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

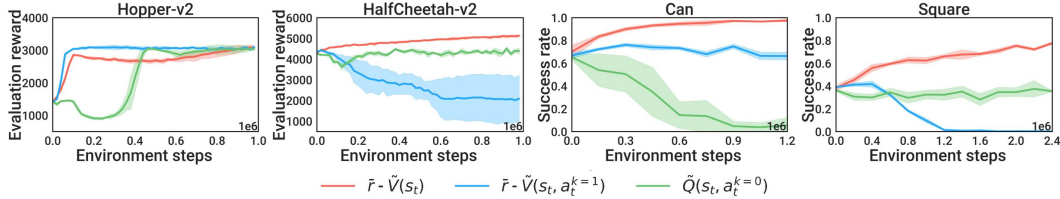
$$\hat{A} = \gamma_{\text{DENOISE}}^k (\bar{r}(\bar{s}_{\bar{t}}, \bar{a}_{\bar{t}}) - \tilde{V}(s_t)).$$

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. (2023)⁶, $\tilde{Q}(s_t, a_t^{k=0})$, that does not directly use the rollout reward \bar{r} in the advantage.

Fig. 13 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

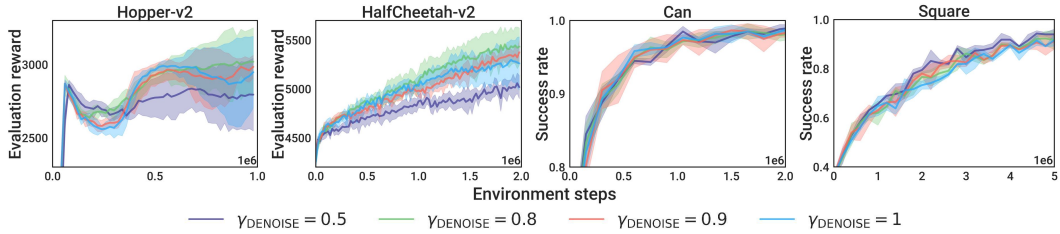
⁶Psenka et al. (2023) applies off-policy training with double Q-learning (according to its open-source implementation) and policy gradient over the denoising steps. Note that this is a baseline in Psenka et al. (2023) that is conjectured to be inefficient. We follow the same except for applying on-policy PPO updates.

1134 **DPPO**'s choice. However, in the more challenging tasks, the environment-state-only advantage
 1135 used in **DPPO** consistently leads to the most improved performance. We believe estimating the
 1136 accurate value of applying a continuous and high-dimensional action can be challenging, and this
 1137 is exacerbated by the high stochasticity of diffusion-based policies and the action chunk size. The
 1138 results here corroborate the findings in Section 5.1 that off-policy Q-learning methods can perform
 1139 well in Hopper-v2 and Walker2D-v2, but exhibit training instability in manipulation tasks from
 1140 ROBOMIMIC.



1141 Figure 13: **Choice of advantage estimator.** Results are averaged over five seeds in Hopper-v2
 1142 and HalfCheetah-v2 and three seeds in Can and Square.

1151 **Denosing discount factor.** We further examine how γ_{DENOISE} in the **DPPO** advantage estimator
 1152 affects fine-tuning. Using a smaller value (i.e., more discount) has the effect of downweighting
 1153 the contribution of earlier denoising steps in the policy gradient. Fig. 14 shows the fine-tuning
 1154 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
 1155 HalfCheetah-v2 $\gamma_{\text{DENOISE}} = 0.8$ leads to better efficiency while smaller $\gamma_{\text{DENOISE}} = 0.5$ slows
 1156 training. The value does not affect training noticeably in Can. In Square the smaller $\gamma_{\text{DENOISE}} =$
 1157 0.5 works slightly better. Overall in manipulation tasks, **DPPO** training seems relatively robust to
 1158 this choice.



1167 Figure 14: **Choice of denoising discount factor.** Results are averaged over five seeds in
 1168 Hopper-v2 and HalfCheetah-v2 and three seeds in Can and Square.

1170 **2. Choice of diffusion noise schedule.** As introduced in Section 4.1, we find it helpful to clip the
 1171 diffusion noise σ_k to a higher minimum value $\sigma_{\text{min}}^{\text{exp}}$ to ensure sufficient exploration. In Figure 15, we
 1172 perform analysis on varying $\sigma_{\text{min}}^{\text{exp}} \in \{.001, .01, .1, .2\}$ (keeping $\sigma_{\text{min}}^{\text{prob}} = .1$ to evaluate likelihoods).
 1173 Although in Can the choice of $\sigma_{\text{min}}^{\text{exp}}$ does not affect the fine-tuning performance, in Square a
 1174 higher $\sigma_{\text{min}}^{\text{exp}} = 0.1$ is required to prevent the policy from collapsing. We conjecture that this is due to
 1175 limited exploration causing policy over-optimizing the collected samples that exhibit limited state-
 1176 action coverage. We also visualize the trajectories at the beginning of fine-tuning in Avoid task
 1177 from D3IL. With higher $\sigma_{\text{min}}^{\text{exp}}$, the trajectories still remain near the two modes of the pre-training
 1178 data but exhibit a higher coverage in the state space — we believe this additional coverage leads to
 1179 better exploration. Anecdotally, we find terminating the denoising process early can also provide
 1180 exploration noise and lead to comparable results, but it requires a more involved implementation
 1181 around the denoising MDP.

1182 **3. Choice of the number of fine-tuned denoising steps.** We examine how the number of fine-
 1183 tuned denoising steps in **DPPO**, K' , affects the fine-tune performance and wall-clock time in Fig. 16.
 1184 We show the curves of individual runs (three for each K') instead of the average as their wall-clock
 1185 times (X-axis) are not perfectly aligned. Generally, fine-tuning too few denoising steps (e.g., 3)
 1186 can lead to subpar asymptotic performance and slower convergence especially in Can. Fine-tuning 10
 1187 steps leads to the overall best efficiency. Similar results are also shown in Fig. 19 with Avoid task.
 Lastly, we note that the GPU memory usage scales linearly with K' .

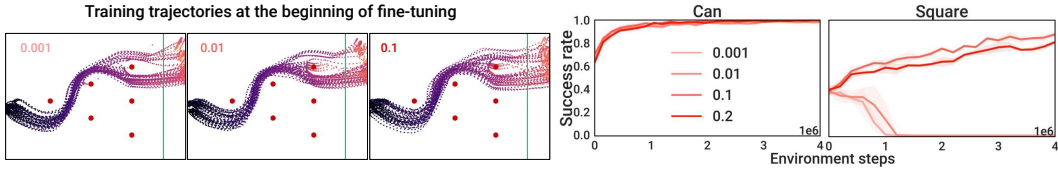


Figure 15: **Choice of minimum diffusion noise.** Results are averaged over three seeds. Note in Left, with higher minimum noise level, the sampled trajectories exhibit wider coverage at the two modes but still maintain the overall structure.

We note that the findings here mostly correlate with those from varying the denoising discount factor, γ_{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 γ_{DENOISE} and K' to achieve the best efficiency.

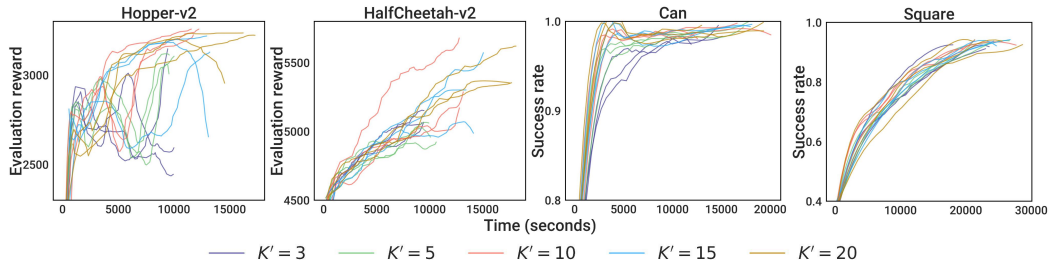


Figure 16: **Choice of number of fine-tuned denoising steps, K' .** Individual runs are shown. The curves are smoothed using a Savitzky–Golay filter.

D.3 EFFECT OF EXPERT DATA

We investigate the effect of the amount of pre-training expert data on fine-tuning performance. In Fig. 17 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.

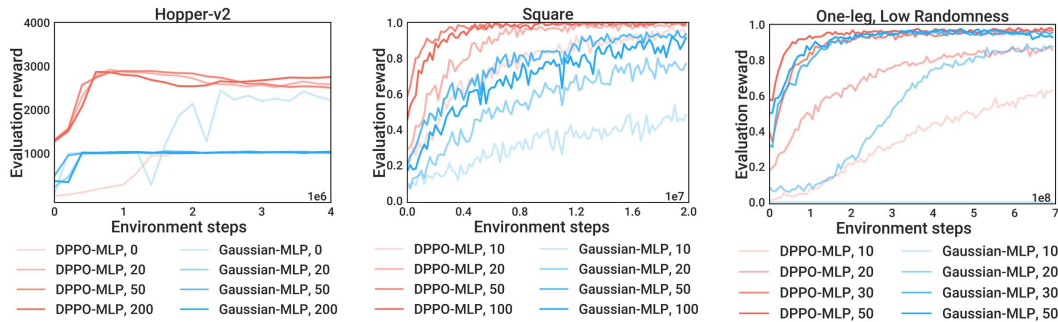


Figure 17: **Varying the number of expert demonstrations.** The numbers in the legends indicates the number of episodes used in pre-training.

Training from scratch. In Fig. 18 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. 17, 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.

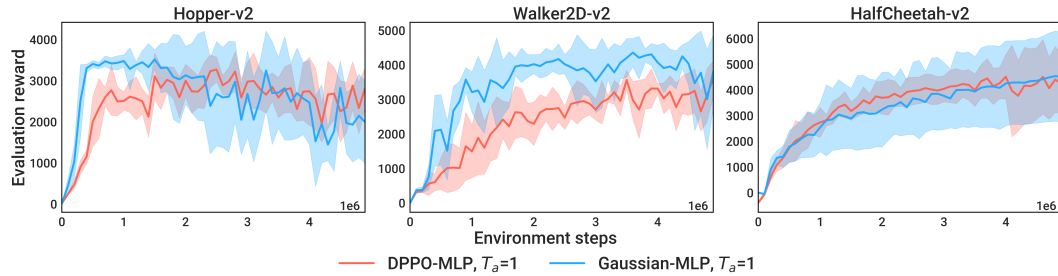


Figure 18: **No expert data / pre-training** with GYM tasks. Results are averaged over five seeds.

D.4 COMPARING TO OTHER POLICY PARAMETERIZATIONS IN `AVOID`

Figure 19 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.

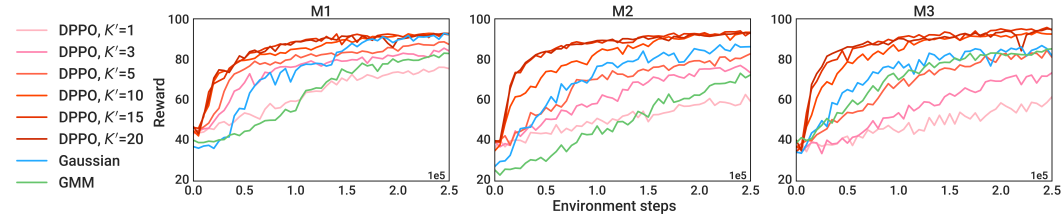


Figure 19: Fine-tuning performance (averaged over five seeds, standard deviation not shown) after pre-training with M1, M2, and M3 in **Avoid task from D3IL**. **DPPO** ($K = 20$), Gaussian, and GMM policies are compared. We also sweep the number of fine-tuned denoising steps K' in **DPPO**.

D.5 COMPARING TO OTHER POLICY PARAMETERIZATIONS IN THE EASIER TASKS FROM ROBOMIMIC

Figure 20 compares the performance of **DPPO** 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 **DPPO** and Gaussian baselines (GMM still exhibits subpar performance). Nonetheless, **DPPO** still achieves similar or even better sample efficiency compared to Gaussian baseline.

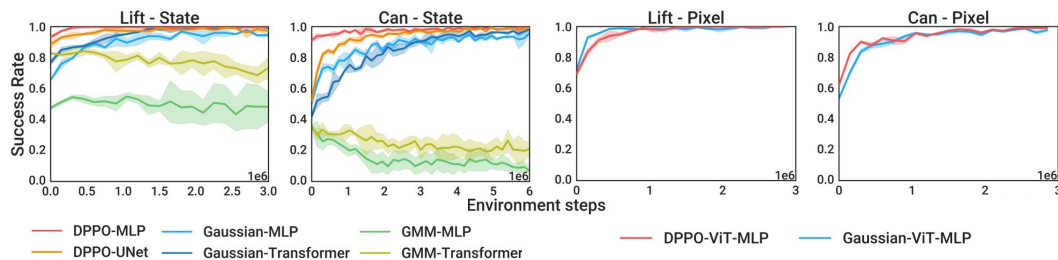


Figure 20: **Comparing to other policy parameterizations** in the easier `Lift` and `Can` tasks from ROBOMIMIC, with **state** (left) or **pixel** (right) observation. Results are averaged over three seeds.

D.6 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 (Song et al., 2020b). We refer the readers to Appendix. D in Song et al. (2020b) for a comprehensive exposition. We follow the official open-source code from Song et al.⁷, and implement policy gradient (single-level MDP) that uses the exact action likelihood $\pi_{\theta}(a_t|s_t)$.

Fig. 21 shows the comparison between **DPPO** and diffusion policy gradient using exact likelihood estimate. Exact policy gradient improves the base policy in `Hopper-v2` but does not outperform **DPPO**. 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 **DPPO** offers, which have shown in Appendix D.2 to often improve fine-tuning efficiency.

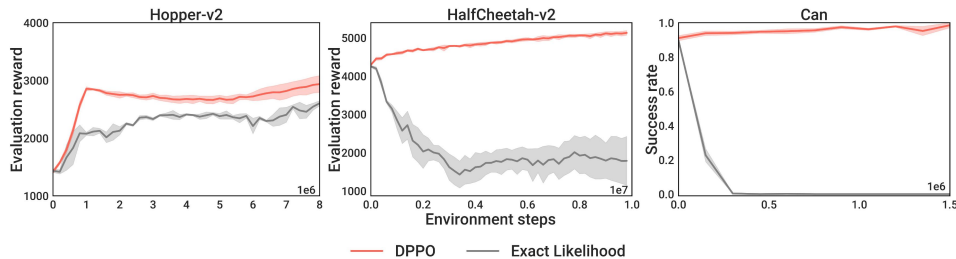


Figure 21: **Comparing to diffusion policy gradient with exact action likelihood.** Results are averaged over five seeds in `Hopper-v2` and `HalfCheetah-v2`, and three seeds in `Can`.

E REPORTING OF WALL-CLOCK TIMES

Comparing to other diffusion-based RL algorithms Section 5.1. Table 1 and Table 2 shows the the wall-clock time used in each OpenAI GYM task and ROBOMIMIC task. In GYM tasks, on average **DPPO** trains 41%, 37%, and 12% faster than **DAWR**, **DIPO**, and **DQL**, respectively, which all require a significant amount of gradient updates per sample to train stably. **QSM**, **DRWR**, and **IDQL** trains 43%, 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. All methods use comparable time except for **DIPO** that uses slightly more on average.

Method	Task		
	Hopper-v2	Walker2D-v2	HalfCheetah-v2
DRWR	11.3	12.7	10.4
DAWR	30.4	30.7	27.1
DIPO	27.8	27.9	26.0
IDQL	16.3	16.1	15.5
DQL	20.5	20.5	17.6
QSM	9.6	9.9	9.7
DPPO	16.6	18.3	16.8

Table 1: **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).

⁷https://github.com/yang-song/score_sde_pytorch

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Method	Task			
	Lift	Can	Square	Transport
DRWR	32.5	39.5	59.8	346.1
DAWR	38.6	46.0	70.5	354.3
DIPO	43.9	51.6	73.3	359.7
IDQL	33.8	41.7	63.7	349.9
DQL	36.9	44.4	68.5	353.5
QSM	31.8	44.5	68.7	322.5
DPPO	35.2	42.0	65.6	350.3

Table 2: **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).

Comparing to other policy parameterizations and architecture Section 5.3 and Section 5.4. Table 3 and Table 4 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 5 shows the wall-clock time used in FURNITURE-BENCH tasks. **DPPO-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
DPPO-MLP	35.6	43.3	65.0	350.5
DPPO-UNet	83.6	92.7	130.4	431.1

Table 3: **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
DPPO-ViT-MLP	194.9	202.5	328.5	871.3

Table 4: **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).

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

Method	Task		
	One-leg	Lamp	Round-table
Gaussian-MLP	101.8	202.8	168.7
DPPO -UNet	148.4	258.2	188.6

Table 5: **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).

F ADDITIONAL EXPERIMENTAL DETAILS

F.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. (2024b), 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 5.3, 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 5.3, we also consider UNet (Ronneberger et al., 2015) as a possible architecture for DP. We follow the implementation from Chi et al. (2024b) 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 5.3 we use Vision-Transformer(ViT)-based image encoder introduced by Hu et al. (2023) before an MLP head. Proprioception input is appended to each channel of the image patches. We also follow (Hu et al., 2023) 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.

F.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 (Todorov et al., 2012). 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 7 and Table 8.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Task	Obs dim - State	Obs dim - Pixel	Act dim	T	Sparse reward ?	
GYM	Hopper-v2	11	-	3	1000	No
	Walker2D-v2	17	-	6	1000	No
	HalfCheetah-v2	17	-	6	1000	No
ROBOMIMIC, state input	Lift	19	-	7	300	Yes
	Can	23	-	7	300	Yes
	Square	23	-	7	400	Yes
	Transport	59	-	14	800	Yes
ROBOMIMIC, pixel input	Lift	9	96×96	7	300	Yes
	Can	9	96×96	7	300	Yes
	Square	9	96×96	7	400	Yes
	Transport	18	2×96×96	14	800	Yes
FURNITURE-BENCH	One-leg	58	-	10	700	Yes
	Lamp	44	-	10	1000	Yes
	Round-table	44	-	10	1000	Yes
D3IL	Avoid	4	-	2	100	Yes

Table 6: **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.

F.3 DESCRIPTIONS OF DIFFUSION-BASED RL ALGORITHM BASELINES IN SECTION 5.1

DRWR: This is a **customized** reward-weighted regression (RWR) algorithm Peters and Schaal (2007) 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 β 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_θ . 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 Peng et al. (2019) that builds on **DRWR** but uses TD-bootstrapped Sutton and Barto (2018) 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 Peng et al. (2019); Black et al. (2023).

The advantage is scaled with β and the exponentiated weight is clipped at w_{\max} . Unlike **DRWR**, we follow (Peng et al., 2019) 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_θ .

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

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

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

where A is calculated using TD(λ), with λ as λ_{DAWR} and the discount factor γ_{ENV} .

DIPO (Yang et al., 2023): 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_{DIPO} iterations with learning rate α_{DIPO} .

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

1512 The actor is then updated with a replay ratio of N_θ .

$$1513 \mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\|\varepsilon_t - \varepsilon_\theta(a_t^{M_{\text{DIPQ}}, k=0}, s_t, k)\|^2].$$

1514 The critic is trained to minimize the Bellman residual with a replay ratio of N_ϕ . Double Q-learning
1515 is also applied.

$$1516 \mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{Q}_\phi(s_{t+1}, \bar{\pi}_\theta(a_{t+1}^{k=0} | s_{t+1})) - \hat{Q}_\phi(s_t, a_t^{m=0, k=0}))\|^2]$$

1517 **IDQL (Hansen-Estruch et al., 2023):** This baseline learns a state-action Q function and state V
1518 function to choose among the sampled actions from DP. DP fits on new samples without weighting.

1519 Again recent experiences are saved in a replay buffer \mathcal{D} . The state value function is updated to match
1520 the expected Q value with an expectile loss, with a replay ratio of N_ψ .

$$1521 \mathcal{L}_\psi = \mathbb{E}^{\mathcal{D}} [|\tau_{\text{IDQL}} - \mathbb{1}(\hat{Q}_\phi(s_t, a_t^0) < \hat{V}_\psi^2(s_t))|].$$

1522 The value function is used to update the Q function with a replay ratio of N_ϕ .

$$1523 \mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{V}_\psi(s_{t+1}) - \hat{Q}_\phi(s_t, a_t^0)\|^2].$$

1524 The actor fits all sampled experiences without weighting, with a replay ratio of N_θ .

$$1525 \mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\|\varepsilon_t - \varepsilon_\theta(a_t^0, s_t, k)\|^2].$$

1526 At inference time, M_{IDQL} actions are sampled from the actor. For training, Boltzmann exploration
1527 is applied based on the difference between Q value of the sampled actions and the V value at
1528 the current state. For evaluation, the greedy action under Q is chosen.

1529 **DQL (Wang et al., 2022):** This baseline learns a state-action Q function and backpropagates the
1530 gradient from the critic through the entire actor (with multiple denoising steps), akin to the usual
1531 Q-learning.

1532 Again recent experiences are saved in a replay buffer \mathcal{D} . The actor is then updated using both a
1533 supervised loss and the value loss with a replay ratio of N_θ .

$$1534 \mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\|\varepsilon_t - \varepsilon_\theta(a_t^0, s_t, k)\|^2 - \alpha_{\text{DQL}} \hat{Q}_\phi(s_t, \bar{\pi}_\theta(a_t^0 | s_t))],$$

1535 where α_{DQL} is a weighting coefficient. The critic is trained to minimize the Bellman residual with a
1536 replay ratio of N_ϕ . Double Q-learning is also applied.

$$1537 \mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{Q}_\phi(s_{t+1}, \bar{\pi}_\theta(a_{t+1}^0 | s_{t+1})) - \hat{Q}_\phi(s_t, a_t^0)\|^2]$$

1538 **QSM (Psenka et al., 2023):** This baseline learns a state-action Q function, and then updates the
1539 actor by aligning the score of the diffusion actor with the gradient of the Q function.

1540 Again recent experiences are saved in a replay buffer \mathcal{D} . The critic is trained to minimize the
1541 Bellman residual with a replay ratio of N_ϕ . Double Q-learning is also applied.

$$1542 \mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{Q}_\phi(s_{t+1}, \bar{\pi}_\theta(a_{t+1}^0 | s_{t+1})) - \hat{Q}_\phi(s_t, a_t^0)\|^2].$$

1543 The actor is updated as follows with a replay ratio of N_θ .

$$1544 \mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [\|\alpha_{\text{QSM}} \nabla_a \hat{Q}_\phi(s_t, a_t) - (-\varepsilon_\theta(a_t^0, s_t, k))\|^2],$$

1545 where α_{QSM} scales the gradient. The negative sign before ε_θ is from taking the gradient of the mean
1546 μ in the denoising process.

1563 F.4 DESCRIPTIONS OF RL FINE-TUNING ALGORITHM BASELINES IN SECTION 5.2

1564 In this subsection, we detail the baselines **RLPD**, **Cal-QL**, and **IBRL**. All policies π_θ are param-
1565 eterized as unimodal Gaussian with an action chunk size of 1.

RLPD (Ball et al., 2023): This baseline is based on Soft Actor Critic (SAC, Haarnoja et al. (2018)) — it learns an entropy-regularized state-action Q function, and then updates the actor by maximizing the Q function w.r.t. the action.

A replay buffer \mathcal{D} is initialized with offline data, and online samples are added to \mathcal{D} . Each gradient update uses a batch of mixed 50/50 offline and online data. An ensemble of N_{critic} critics is used, and at each gradient step two critics are randomly chosen. The critics are trained to minimize the Bellman residual with replay ratio N_ϕ :

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{Q}_{\phi'}(s_{t+1}, \pi_\theta(a_{t+1}|s_{t+1})) - \hat{Q}_\phi(s_t, a_t)\|^2].$$

The target critic parameter ϕ' is updated with delay. The actor minimizes the following loss with a replay ratio of 1:

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [-\hat{Q}_\phi(s_t, a_t) + \alpha_{\text{ent}} \log \pi_\theta(a_t|s_t)],$$

where α_{ent} is the entropy coefficient (automatically tuned as in SAC starting at 1).

Cal-QL (Nakamoto et al., 2024): This baseline trains the policy μ and the action-value function Q^μ in an offline phase and then an online phase. During offline phase only offline data is sampled for gradient update, while during online phase mixed 50/50 offline and online data are sampled. The critic is trained to minimize the following loss (Bellman residual and calibrated Q-learning):

$$\begin{aligned} \mathcal{L}_\phi = & \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \hat{Q}_{\phi'}(s_{t+1}, \pi_\theta(a_{t+1}|s_{t+1}))) - \hat{Q}_\phi(s_t, a_t)\|^2] \\ & + \beta_{\text{cql}} (\mathbb{E}^{\mathcal{D}} [\max(Q_\phi(s_t, a_t), V(s_t))] - \mathbb{E}^{\mathcal{D}} [Q_\phi(s_t, a_t)]), \end{aligned}$$

where β_{cql} is a weighting coefficient between Bellman residual and calibration Q-learning and $V(s_t)$ is estimated using Monte-Carlo returns. The target critic parameter ϕ' is updated with delay. The actor minimizes the following loss:

$$\mathcal{L}_\theta = \mathbb{E}^{\mathcal{D}} [-\hat{Q}_\phi(s_t, a_t) + \alpha_{\text{ent}} \log \pi_\theta(a_t|s_t)],$$

where α_{ent} is the entropy coefficient (automatically tuned as in SAC starting at 1).

IBRL (Hu et al., 2023): This baseline first pre-trains a policy μ_ψ using behavior cloning, and for fine-tuning it trains a RL policy π_θ initialized as μ_ψ . During fine-tuning recent experiences are saved in a replay buffer \mathcal{D} . An ensemble of N_{critic} critics is used, and at each gradient step two critics are randomly chosen. The critics are trained to minimize the Bellman residual with replay ratio N_ϕ :

$$\mathcal{L}_\phi = \mathbb{E}^{\mathcal{D}} [\|(R_t + \gamma_{\text{ENV}} \max_{a' \in \{a^{IL}, a^{RL}\}} \hat{Q}_{\phi'}(s_{t+1}, a') - \hat{Q}_\phi(s_t, a_t)\|^2]$$

where $a^{IL} = \mu_\psi(s_{t+1})$ (no noise) and $a^{RL} \sim \pi_{\theta'}(s_{t+1})$, and $\pi_{\theta'}$ is the target actor. The target critic parameter ϕ' is updated with delay. The actor minimizes the following loss with a replay ratio of 1:

$$\mathcal{L}_\theta = -\mathbb{E}^{\mathcal{D}} [\hat{Q}_\phi(s_t, a_t)].$$

The target actor parameter θ' is also updated with delay.

F.5 ADDITIONAL DETAILS OF **DPPO** IMPLEMENTATION IN ALL TASKS

Similar to all baselines in Appendix F.3, we denote N_θ and N_ϕ 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 (Huang et al., 2022), the batch updates in an iteration terminate when the KL divergence between π_θ and $\pi_{\theta_{\text{old}}}$ reaches 1.

We also find the PPO clipping ratio, ε , 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 ε in $\{.1, .01, .001\}$ and choose the highest value that satisfies the clipping fraction target. Empirically we also find that, using a higher ε for earlier denoising steps in **DPPO** further improves training stability in manipulation tasks. Denote ε_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 .

F.6 ADDITIONAL DETAILS OF ROBOMIMIC TASKS AND TRAINING IN SECTION 5.3

Tasks. We consider four tasks from the ROBOMIMIC benchmark (Mandlekar et al., 2021): (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 (Todorov et al., 2012). Success rate curves shown in Fig. 5, Fig. 6, and Fig. 20 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 6). Detailed hyperparameters are listed in Table 10.

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 11.

F.7 DESCRIPTIONS OF POLICY PARAMETERIZATION BASELINES IN SECTION 5.3

Gaussian. We consider unimodal Gaussian with diagonal covariance, the most commonly used policy parameterization in RL. The standard deviation for each action dimension, σ_{Gau} , is fixed during pre-training; we also tried to learn σ_{Gau} from the dataset but we find the training very unstable. During fine-tuning σ_{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 discusses in Appendix F.5, we choose the PPO clipping ratio ε 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_{GMM} the number of mixtures. The standard deviation for each action dimension in each mixture, σ_{GMM} , is also fixed during pre-training. Again during fine-tuning σ_{GMM} is learned starting from the same fixed value and also clipped between 0.01 and 0.2.

F.8 ADDITIONAL DETAILS OF FURNITURE-BENCH TASKS AND TRAINING IN SECTION 5.4

Tasks. We consider three tasks from the FURNITURE-BENCH benchmark (Heo et al., 2023): (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. 22 for the visualized rollouts in simulation.

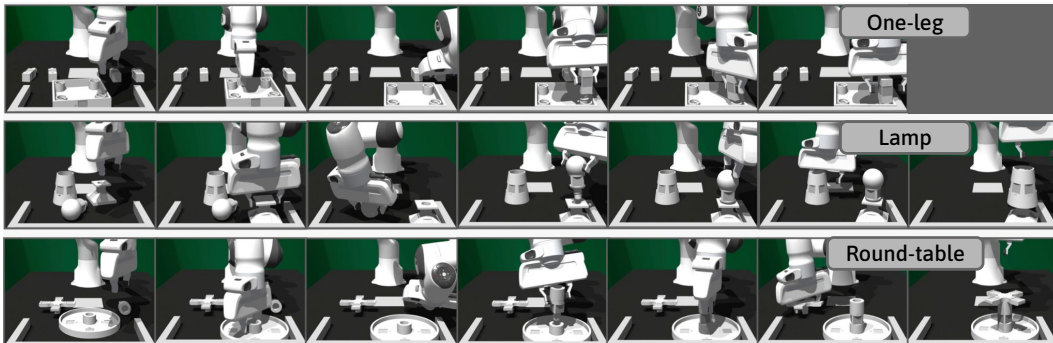
Pre-training. The pre-training dataset is collected in the simulated environments using a Space-Mouse⁸, a 6 DoF input device. The simulator runs at 10Hz. At every timestep, we read off the state

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

1674 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
 1675 before passed to the environment step and stored as the action alongside the current observation in
 1676 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
 1677 not record any action nor pass it to the environment. Discarding no-ops is important for allowing
 1678 the policies to learn from demonstrations effectively. When the desired number of demonstrations
 1679 has been collected (typically 50), we process the actions to convert the delta actions stored from the
 1680 SpaceMouse into absolute pose actions by applying the delta action to the current EE pose at each
 1681 timestep.

1682 The observations and actions are normalized to $[-1, 1]$ using min/max statistics from the pre-training
 1683 dataset. No history observation (proprioception or ground-truth object states) is used, i.e., only the
 1684 current observation is passed to the policy. All policies are trained with batch size 256, learning rate
 1685 $1e-4$ decayed to $1e-5$ with a cosine schedule, and weight decay $1e-6$. Diffusion-based policies are
 1686 trained with 8000 epochs, while Gaussian policies are trained with 3000 epochs. Gaussian policies
 1687 can easily overfit the pre-trained dataset, while diffusion-based policies are more resilient. Gaussian
 1688 policies also require a very large MLP (~ 10 million parameters) to fit the data well.

1689 **Fine-tuning.** Diffusion-based and Gaussian pre-trained policies are then fine-tuned using online
 1690 experiences sampled from 1000 parallelized IsaacGym environments Makoviychuk et al. (2021).
 1691 Success rate curves shown in Fig. 7 are evaluated by running fine-tuned policies with $\sigma_{\min}^{\text{exp}} = 0.001$
 1692 (i.e., without extra noise) for 1000 episodes. Episodes terminate only when they reach maximum
 1693 episode length (shown in Table 6). Detailed hyperparameters are listed in Table 12. We find a smaller
 1694 amount of exploration noise (we set $\sigma_{\min}^{\text{exp}}$ and σ_{Gau} to be 0.04) is necessary for the pre-trained policy
 1695 achieving nonzero success rates at the beginning of fine-tuning.



1709 Figure 22: Representative rollouts from simulated FURNITURE-BENCH tasks.

1710
1711
1712 **Hardware setup - robot control.** The physical robot used is a Franka Emika Panda arm. The
 1713 policies output a sequence of desired end-effector poses in the robot base frame to control the robot.
 1714 These poses are converted into joint position targets through differential inverse kinematics. We
 1715 calculate the desired end-effector velocity as the difference between the desired and current poses
 1716 divided by the delta time $dt = 1/10$. We then convert this to desired joint velocities using the
 1717 Jacobian and compute the desired joint positions with a first-order integration over the current joint
 1718 positions and desired velocity. The resulting joint position targets are passed to a low-level joint
 1719 impedance controller provided by Polymetis (Lin et al., 2021), running at 1kHz.

1720 **Hardware setup - state estimation.** To deploy state-based policies on real hardware, we utilize
 1721 AprilTags (Wang and Olson, 2016) for part pose estimation. The FURNITURE-BENCH (Heo et al.,
 1722 2023) task suite provides AprilTags for each part and code for estimating part poses from tag detec-
 1723 tions. The process involves several steps: (1) detecting tags in the camera frame, (2) mapping tag
 1724 detections to the robot frame for policy compatibility, (3) utilizing known offsets between tags and
 1725 object centers in the simulator, and (4) calibrating the camera pose using an AprilTag at a known
 1726 position relative to the robot base. Despite general accuracy, detections can be noisy, especially
 1727 during movement or partial occlusion, which the `One-leg` task features. Since the task requires
 high precision, we find the following to help make the estimation reliable enough:

-
- 1728
- 1729
- 1730
- 1731
- 1732
- 1733
- 1734
- 1735
- 1736
- 1737
- 1738
- 1739
- 1740
- 1741
- 1742
- 1743
- 1744
- 1745
- 1746
- 1747
- 1748
- 1749
- 1750
- 1751
- 1752
- 1753
- 1754
- 1755
- 1756
- 1757
- 1758
- 1759
- 1760
- 1761
- 1762
- 1763
- 1764
- 1765
- 1766
- 1767
- 1768
- 1769
- 1770
- 1771
- 1772
- 1773
- 1774
- 1775
- **Camera coverage:** We find detection quality sensitive to distance and angle between the camera and tag. This issue is likely due to the RealSense D435 camera having mediocre image quality and clarity and the relatively small tags. To remedy this, we opt to use 4 cameras roughly evenly spread out around the scene to ensure that at least one camera has a solid view of a tag on all the parts (i.e., as close as possible with a straight-on view). To find the best camera positions, we start with having a camera in each of the cardinal directions around the scene. Then, we adjust the pose of each to get it as close as possible to the objects while still covering the necessary workspace and capturing the base tag for calibration. Moving the robot arm around the scene to avoid the worst occlusion is also helpful.
 - **Lighting:** Even with better camera coverage and placement, detection quality depends on having crisp images. We find proper lighting helpful to improve image quality. In particular, the scene should be well and evenly lit around the scene without causing reflections in either the tag or table.
 - **Filtering:** Bad detections can sometimes cause the resulting pose estimate to deviate significantly from the true pose, i.e., jumping several centimeters from one frame to the next. This usually only happens on isolated frames, and thus before “accepting” a given detection, we check if the new position and orientation are within 5 cm and 20 degrees of the previously accepted pose. In addition, we apply low-pass filtering on the detection using a simple exponential average (with $\alpha = 0.25$) to smooth out the high-frequency noise.
 - **Averaging:** The objects have multiple tags that can be detected from multiple cameras. After performing the filtering step, we average all pose estimates for the same object across different tags and cameras, which also helps smooth out noise. This alone, however, does not fully cancel the case when a single detection has a large jump, as this can severely skew the average, still necessitating a filtering step. Having multiple cameras benefits this step, too, as it provides more detections to average over.
 - **Caching part pose in hand:** A particularly difficult phase of the task to achieve good detections is when the robot transports the table leg from the initial position to the tabletop for insertion. The main problems are that the movement can blur the images, and the grasping can cause occlusions. Therefore, we found it helpful to assume that once the part was grasped by the robot, it would not move in the grasp until the gripper opened. With this, we can “cache” the pose of the part relative to the end-effector once the object is fully grasped and use this instead of relying on detections during the movement.
 - **Normalization pitfalls and clipping:** We generally use min-max normalization of the state observations to ensure observations are in $[-1, 1]$. The tabletop part moves very little in the z -direction demonstration data, meaning the resulting normalization limits (the minimum and maximum value of the data) can be very close, $x_{\max} - x_{\min} \approx 0$. With these tight limits, the noise in the real-world detection can be amplified greatly as $x_{\text{norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$. Therefore, ensure that normalization ranges are reasonable. As an extra safeguard, clipping the data to $[-1, 1]$ can also help.
 - **Only estimate necessary states:** Despite the `One-leg` task having 5 parts, only 2 are manipulated. Only estimating the pose of those parts can eliminate a lot of noise. In particular, the pose of the 3 legs that are not used and the obstacle (the U-shaped fixture) can be set to an arbitrary value from the dataset.
 - **Visualization for debugging:** We use the visualization tool MeshCat⁹ extensively for debugging of state estimation. The tool allows for easy visualizations of poses of all relevant objects in the scene, like the robot end-effector and parts, which makes sanity-checking the implementation far easier than looking at raw numbers.

1776 **Hardware evaluation.** We perform 20 trials for each method. We adopt a single-blind model
1777 selection process: at the beginning of each trial, we first randomize the initial state. Then, we
1778 randomly select a method and roll it out, but the experimenter does not observe which model is
1779 used. We record the success and failure of each trial and then aggregate statistics for each model
1780 after all trials are completed.

1781

⁹<https://github.com/meshcat-dev/meshcat>

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835



Figure 23: **Qualitative comparison of pre-trained vs. fine-tuned DPPO policies in real 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 that requires corrective behavior.

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 **DPPO** to simulate the imperfect low-level controller; we find adding such noise to the Gaussian policy leads to zero task success rate while **DPPO** 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 behavior cloning (BC) regularization loss in fine-tuning with the Gaussian baseline. The combined loss follows

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

where $\pi_{\theta_{pre-trained}}$ 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 $\alpha_{BC} = 0.1$. 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.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

F.9 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 (Todorov et al., 2012). Reward curves shown in Fig. 9 and Fig. 19 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. 9 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.

F.10 LISTED TRAINING HYPERPARAMETERS

Method	Parameter	Task(s)			
		GYM	Lift, Can	Square	Transport
Common	γ_{ENV}	0.99	0.999	0.999	0.999
	σ_{min}^{exp}	0.1	0.1	0.1	0.08
	σ_{min}^{prob}			0.1	
	T_a	4	4	4	8
	K			20	
	Actor learning rate	1e-4	1e-5	1e-5	1e-5 (decayed to 1e-6)
	Critic learning rate (if applies)			1e-3	
	Actor MLP dims	[512, 512, 512]	[512, 512, 512]	[1024, 1024, 1024]	[1024, 1024, 1024]
Critic MLP dims (if applies)			[256, 256, 256]		
DRWR	β			10	
	w_{max}			100	
	N_θ			16	
	Batch size			1000	
DAWR	β			10	
	w_{max}			100	
	λ_{DAWR}			0.95	
	N_θ			64	
	N_ϕ	16	4	4	4
	Buffer size	200000	150000	150000	150000
DIPO	Batch size			256	
	α_{DIPO}			1e-4	
	M_{DIPO}			10	
	N_θ			64	
IDQL	Buffer size			400000	
	Batch size			5000	
	M_{IDQL}	20	10	10	10
	N_θ			16	
	N_ϕ			16	
DQL	Buffer size	200000	150000	150000	150000
	Batch size	256	512	512	512
	α_{DQL}			1	
	N_θ			64	
DQL	N_ϕ			64	
	Buffer size			400000	
	Batch size			5000	

Table 7: 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
QSM	α_{QSM}			50	
	N_θ			32	
	N_ϕ			32	
	Buffer size	200000	150000	150000	150000
	Batch size			5000	
DPPO	$\gamma_{DENOISE}$			0.99	
	GAE λ			0.95	
	N_θ	5	10	10	8
	N_ϕ	5	10	10	8
	ε			0.01	
	Batch size	5000	7500	10000	10000
	K'			10	

Table 8: Continuation of Table 7.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

Method	Parameter	Task(s)		
		HalfCheetah-v2	Can	Square
Common	γ_{ENV}	0.99	0.999	0.999
	T_a		1	
RLPD	N_ϕ	20	3	3
	N_{critic}	10	5	5
	Batch size		256	
Cal-QL	β_{cql}		5	
	Batch size		256	
IBRL	N_ϕ	5	3	3
	N_{critic}		5	
	Batch size		256	
DPPO	σ_{min}^{exp}		0.1	
	σ_{min}^{prob}		0.1	
	$\gamma_{DENOISE}$		0.99	
	GAE λ		0.95	
	N_θ	5	10	10
	N_ϕ	5	10	10
	ϵ		0.01	
	Batch size	5000	7500	10000
	K		20	
	K'		10	

Table 9: Fine-tuning hyperparameters for HalfCheetah-v2, Can, and Square when **comparing demo-augmented RL methods**. We list hyperparameters shared by all methods first, and then method-specific ones.

Method	Parameter	Task		
		Lift, Can	Square	Transport
Common	γ_{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 λ		0.95	
	N_θ	10	10	8
	N_ϕ	10	10	8
	ϵ		0.01 (annealed in DPPO)	
	Batch size	7500	10000	10000
Gaussian, Common	σ_{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_{GMM}		5	
	σ_{GMM}	0.1	0.1	0.08
	Model size	1.15M	4.40M	4.90M
GMM-Transformer	Model size	680K	1.87M	1.89M
DPPO, Common	$\gamma_{DENOISE}$		0.99	
	σ_{min}^{exp}	0.1	0.1	0.08
	σ_{min}^{prob}	0.1	0.1	0.1
	K		20	
	K'		10	
DPPO-MLP	Model size	576K	2.31M	2.43M
DPPO-UNet	Model size	652K	1.62M	1.68M

Table 10: 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.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

Method	Parameter	Task		
		Lift, Can	Square	Transport
Common	γ^{ENV}		0.999	
	T_α	4	4	8
	Actor learning rate	1e-4	1e-5	1e-5 (decayed to 1e-6)
	Critic learning rate			1e-3
	GAE λ			0.95
	N_θ	10	10	8
	N_ϕ	10	10	8
	ε		0.01 (annealed in DPPO)	
Gaussian-ViT-MLP	Batch size	7500	10000	10000
	Model size	1.03M	1.03M	1.93M
Gaussian-ViT-MLP	σ_{Gau}	0.1	0.1	0.08
	Model size	1.06M	1.06M	2.05M
DPPO -ViT-MLP	γ^{DENOISE}		0.9	
	$\sigma_{\text{min}}^{\text{exp}}$	0.1	0.1	0.08
	$\sigma_{\text{min}}^{\text{prob}}$		0.10	
	K		100	
	K'		5 (DDIM)	

Table 11: 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	γ^{ENV}		0.999	
	T_α		8	
	Actor learning rate		1e-5 (decayed to 1e-6)	
	Critic learning rate			1e-3
	GAE λ			0.95
	N_θ			5
	N_ϕ			5
	ε		0.001	
Gaussian-MLP	Batch size		8800	
	Model size	10.64M	10.62M	10.62M
Gaussian-MLP	σ_{Gau}		0.04	
	Model size	6.86M	6.81M	6.81M
DPPO -UNet	γ^{DENOISE}		0.9	
	$\sigma_{\text{min}}^{\text{exp}}$		0.04	
	$\sigma_{\text{min}}^{\text{prob}}$		0.1	
	K		100	
	K'		5 (DDIM)	

Table 12: Fine-tuning hyperparameters for FURNITURE-BENCH tasks when **comparing policy parameterizations**. We list hyperparameters shared by all methods first, and then method-specific ones.