

LATENT WEIGHT DIFFUSION: GENERATING POLICIES FROM TRAJECTORIES

Anonymous authors

Paper under double-blind review

ABSTRACT

With the increasing availability of open-source robotic data, imitation learning has emerged as a viable approach for both robot manipulation and locomotion. Currently, large generalized policies are trained to predict controls or trajectories using diffusion models, which have the desirable property of learning multimodal action distributions. However, generalizability comes with a cost — namely, larger model size and slower inference. Further, there is a known trade-off between performance and action horizon for Diffusion Policy (i.e., diffusing trajectories): fewer diffusion queries accumulate greater trajectory tracking errors. Thus, it is common practice to run these models at high inference frequency, subject to robot computational constraints.

To address these limitations, we propose Latent Weight Diffusion (LWD), a method that uses diffusion to learn a distribution over policies for robotic tasks, rather than over trajectories. Our approach encodes demonstration trajectories into a latent space and then decodes them into policies using a hypernetwork. We employ a diffusion denoising model within this latent space to learn its distribution. We demonstrate that LWD can reconstruct the behaviors of the original policies that generated the trajectory dataset. LWD offers the benefits of considerably smaller policy networks during inference and requires fewer diffusion model queries. When tested on the Metaworld MT10 benchmark, LWD achieves a higher success rate compared to a vanilla multi-task policy, while using models up to $\sim 18x$ smaller during inference. Additionally, since LWD generates closed-loop policies, we show that it outperforms Diffusion Policy in long action horizon settings, with reduced diffusion queries during rollout.

1 INTRODUCTION

The recent increase in open-source robotic data has made imitation learning an attractive prospect to solve robot manipulation and locomotion tasks (Collaboration et al., 2023; Peng et al., 2020). While traditional supervised learning methods like Behavioral Cloning (Florence et al., 2022) and transformer-based models such as RT-1 (Brohan et al., 2022) have demonstrated some success, they are unable to capture the multimodal nature of robotic action distributions (e.g., while avoiding an obstacle in a navigation task with two optimal actions in opposing directions ‘turn left’ and ‘turn right’, the learned action ‘go straight’ is a suboptimal average of the two). Recently, diffusion-based methods have emerged as a promising alternative for robot control (Tan et al., 2024), offering the advantages of continuous outputs and the capacity to learn multimodal action distributions.

Inspired by the success of latent diffusion in vision (Rombach et al., 2022b) and language (Lovelace et al., 2024), we explore it here for robotics. We introduce a novel method that utilizes diffusion models to learn a distribution of policies for robotic tasks from demonstration data. Unlike existing robotics approaches focusing on trajectory diffusion (Chi et al., 2024), our method Latent Weight Diffusion (LWD) diffuses neural network weights to generate policies. This is achieved by encoding demonstration trajectories into a latent space, employing a diffusion denoising model to learn the distribution of latents within this space, followed by decoding the latent representations into executable policies using a hypernetwork (Ha et al., 2016).

LWD provides several benefits. First, the distribution modeling capabilities of diffusion models allow it to capture complex behavior distributions. Second, generating the parameters of a neu-

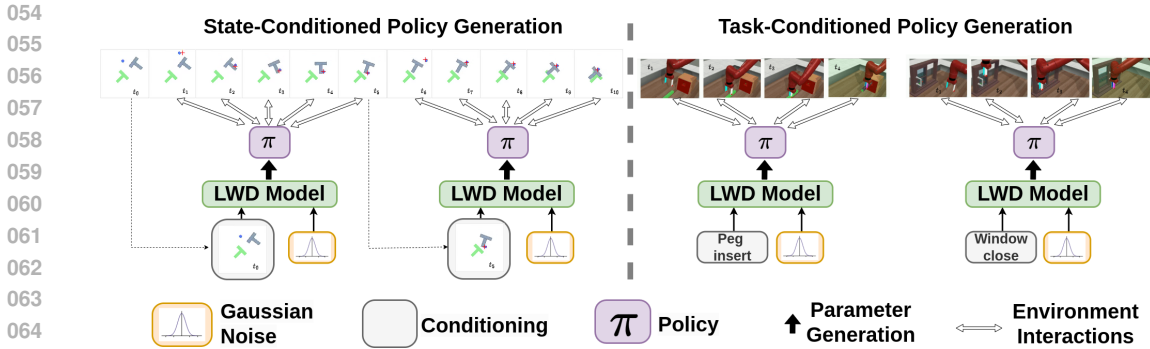


Figure 1: **Latent Weight Diffusion (LWD)** generates policies from heterogeneous trajectory data. With state-conditioned policy generation, the diffusion model can run inference at a lower frequency. With task-conditioned policy generation, the generated policies can be small yet maintain task-specific performance. Demonstrations of this work can be found on the project website: <https://sites.google.com/view/iclr2024submission/home>

ral network policy enables them to be run as closed-loop controllers. Since closed-loop control is less susceptible to trajectory tracking errors than open-loop control, LWD allows for longer action horizons than trajectory generation methods. This has the additional advantage of reducing calls to an expensive diffusion model. Third, the LWD model can be conditioned on task identifiers to generate task-specific policies. This provides a unique advantage of maintaining multi-task generalization in the diffusion model and task-specific performance in the generated policy. Since the generated policy is task-specific, it has fewer parameters than a generalist multi-task agent during inference. These advantages demonstrate the effectiveness of latent diffusion in learning policies from trajectory data, leading to the following performance gains.

1. **Policy Diversity:** LWD learns a diverse policy distribution from a trajectory dataset. On the D4RL benchmark, when trained on a trajectory mixture from three behaviors, LWD accurately generates policies capturing the behavior distributions of the original policies.
2. **Closed-loop Control:** LWD generates closed-loop policies enabling reactivity to environmental changes. This enables policies generated by LWD to run for 4X longer action horizons, compared to a Diffusion Policy model for the same change in performance, when evaluated on the PushT task.
3. **Small Policy Size:** The burden of generalization is borne by the diffusion model instead of the inference policy, allowing the generated policy to have fewer parameters, without compromising multi-task performance. On the Metaworld MT10 suite of tasks, LWD shows similar multi-task performance as a BC baseline that is up to 18X larger.

2 RELATED WORK

2.1 IMITATION LEARNING FOR ROBOTICS

With the availability of vast swaths of open-source robotic data, imitation learning has emerged as a viable approach for robot control. Original imitation learning approaches were simple – behavioral cloning agents that learned to predict controls or trajectories from expert demonstrations. With the advent of transformer-based methods, imitation learning has grown in popularity. Methods like PerAct (Shridhar et al., 2022) and RT-1 (Brohan et al., 2022) perform well on various tasks. Brohan et al. (2023) showed that VLMs can be combined with robot demonstrations to solve tasks by directly interpreting the token output as actions. Collaboration et al. (2023) showed that transformer-based methods can be used to learn from demonstrations across different embodiments. Object-aware policies like Heravi et al. (2022) have been shown to improve performance for visuomotor tasks.

2.2 DIFFUSION

Diffusion models have emerged as a leading approach in image generation, with Denoising Diffusion Probabilistic Models (DDPM) being a prominent class of generative models (Ho et al., 2020).

108 These models generate images by progressively denoising samples drawn from an isotropic Gaus-
 109 sian distribution. Furthermore, Rombach et al. (2022a) demonstrated that diffusion can be effec-
 110 tively applied in the latent space of a pre-trained Variational Autoencoder. Recently, diffusion-based
 111 methods have shown promise in solving robotic tasks. The seminal work Chi et al. (2024) showed
 112 that diffusion models can be used to learn multimodal action distributions for a task by diffusing tra-
 113 jectories for control up to a defined action horizon. Urain et al. (2022) showed that diffusion models
 114 can be used to learn smooth cost functions for the joint optimization of grasp and motion plans.
 115 For long horizon skills, Mishra et al. (2023) showed how to use diffusion to chain skills together to
 116 solve a larger task. Diffusion models have been used to formulate policies to control a quadruped
 117 robot Huang et al. (2024), although the length of the trajectories diffused was still relatively short,
 118 therefore requiring a higher diffusion inference frequency. Tan et al. (2024) showed that latent dif-
 119 fusion could be applied to multi-task manipulation action trajectory generation. A key limitation of
 120 using diffusion models to generate trajectories is the degradation in performance for longer action
 121 horizons, which is caused by both modeling and trajectory tracking errors.

122 2.3 HYPERNETWORKS

124 Hypernetworks, introduced by Ha et al. (2016), are neural networks that can estimate the weights
 125 of a secondary network. Following their inception, they have been extended and applied in multiple
 126 settings. In meta-learning, Bertinetto et al. (2016) proposed a model where a learner network pre-
 127 dicts the parameters of another network for one-shot learning tasks, sharing conceptual similarities
 128 with hypernetworks. The concept of dynamically generating network parameters is also related to
 129 Dynamic Filter Networks by Jia et al. (2016), where filters are generated on the fly based on the
 130 input. This method aligns with the principles of hypernetworks, emphasizing adaptability and ef-
 131 ficiency in processing varying inputs. It has also been shown that hypernetworks can be used for
 132 robot policy representation (Hegde et al., 2024).

133 2.4 MULTI-TASK LEARNING

135 Metaworld (Yu et al., 2020) and RLBench (James et al., 2019) are popular benchmarks for multi-task
 136 learning. Many recent transformer-based methods have shown good performance on various tasks,
 137 given task conditioning during training and testing, such as Shridhar et al. (2022) and the RT family
 138 of models. GNFactor (Ze et al., 2023) uses a generalizable neural feature field to learn a volumetric
 139 representation of the environment, which can be used to synthesize different views of the environ-
 140 ment. Another way to solve multi-task learning is through modularity. Devin et al. (2016) showed
 141 how to split networks into modules that are task-specific and robot-specific. Naturally, the task-
 142 specific modules can be shared across robots, and the robot-specific modules can be shared across
 143 tasks on a robot, enabling the transfer of learned behaviors across tasks or robot embodiments.

144 3 PROBLEM FORMULATION & METHOD

147 This work builds on Hegde et al. (2023), which demonstrates the capacity of latent diffusion mod-
 148 els to generate policies from a policy dataset while addressing its key limitation - the reliance on
 149 often unavailable policy datasets - by utilizing trajectory datasets instead. LWD employs a two-
 150 step process. A variational autoencoder (VAE) with a weak KL-regularization coefficient encodes
 151 trajectories into a latent space that can be decoded into a trajectory. A diffusion model learns the dis-
 152 tribution of this latent space, enabling policy sampling from the learned distribution (see Figure 2).

153 3.1 LATENT POLICY REPRESENTATION

155 Consider $\pi(\cdot, \theta)$ as a stochastic policy, parameterized by θ , that interacts with the environment and
 156 generates trajectories τ . Suppose there exists a distribution of policy parameters $p(\theta)$, and we sample
 157 a policy parameter from this distribution and collect a trajectory for this sampled policy parameter.
 158 This process of sampling the policy parameter is done for all trajectories collected. We assume that
 159 for a given θ , $a_t = \pi(s_t, \theta) + e$, where e is normally distributed around 0. i.e., $a_t \sim \mathcal{N}(\pi(s_t, \theta), \sigma^2)$
 160

161 Our goal is to learn the distribution of policy parameters $p(\theta)$ that produced the trajectory dataset.
 We assume that a latent variable z exists that contains information required to identify different

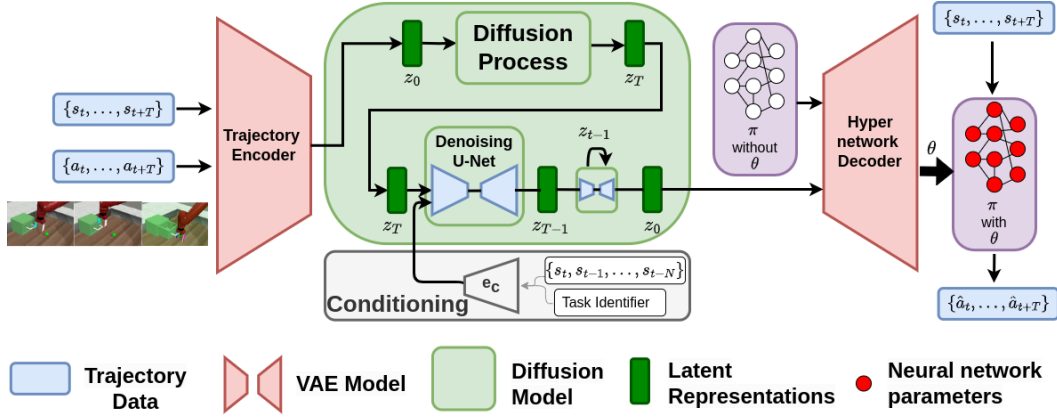


Figure 2: **LWD**: We first pre-train a VAE that variationally encodes trajectories to a latent space and then decodes it as policy parameters. Next, we train a conditional latent diffusion model to learn this latent distribution.

policy behaviors. Since trajectories are generated using parameters θ , we can use conditional independence, $p(\tau | z, \theta) = p(\tau | \theta)$. Considering that our dataset consists of trajectories, we want to maximize the likelihood of sampling τ , therefore maximizing $\log p(\tau)$. We derive a modified version of the Evidence Lower Bound (ELBO) to incorporate $p(\theta)$. This is shown below.

$$\begin{aligned}
 \log p(\tau) &= \log \int \int p(\tau, \theta, z) dz d\theta \quad (\text{Introduce policy parameter } \theta \text{ and latent variable } z) \\
 &= \log \int \int p(\tau | z, \theta) p(\theta | z) p(z) dz d\theta \quad (\text{Apply the chain rule}) \\
 &= \log \int \int \frac{p(\tau | z, \theta) p(\theta | z) p(z)}{q(z | \tau)} q(z | \tau) dz d\theta \quad (1a) \\
 &\quad (\text{Introduce a variational distribution } q(z | \tau), \text{ approximating the true posterior } p(z | \tau)) \\
 &= \log \int \mathbb{E}_{p(\theta|z)} \left[\frac{p(\tau | z, \theta) p(z)}{q(z | \tau)} q(z | \tau) \right] dz \quad (1b) \\
 &= \log \mathbb{E}_{q(z|\tau)} \left[\frac{\mathbb{E}_{p(\theta|z)} [p(\tau | z, \theta)] p(z)}{q(z | \tau)} \right] \quad (1c) \\
 &\geq \mathbb{E}_{q(z|\tau)} \left[\log \left(\frac{\mathbb{E}_{p(\theta|z)} [p(\tau | z, \theta)] p(z)}{q(z | \tau)} \right) \right] \quad (\text{Jensen's inequality}) \\
 &= \mathbb{E}_{q(z|\tau)} [\log (\mathbb{E}_{p(\theta|z)} [p(\tau | z, \theta)])] - \mathbb{E}_{q(z|\tau)} [\log (q(z | \tau)) - \log (p(z))] \quad (1d) \\
 &= \mathbb{E}_{q(z|\tau)} [\log (\mathbb{E}_{p(\theta|z)} [p(\tau | \theta)])] - \text{KL}(q(z | \tau) \| p(z)) \quad (\text{conditional independence}) \\
 &\quad (1e) \\
 &\geq \mathbb{E}_{q(z|\tau)} [\mathbb{E}_{p(\theta|z)} [\log (p(\tau | \theta))]] - \text{KL}(q(z | \tau) \| p(z)) \quad (\text{Jensen's inequality}) \\
 &\quad (1f)
 \end{aligned}$$

Assuming the state transitions are Markov and s_1 is independent of θ , the joint likelihood of the entire sequence $\{(s_1, a_1), (s_2, a_2), \dots, (s_T, a_T)\}$ (i.e., $p(\tau | \theta)$) is given by:

$$p(s_1, a_1, \dots, s_T, a_T | \theta) = p(s_1) p(a_1 | s_1, \theta) \cdot \prod_{t=2}^T p(s_t | s_{t-1}, a_{t-1}) p(a_t | s_t, \theta) \quad (2a)$$

$$\log p(s_1, a_1, \dots, s_T, a_T | \theta) = \log p(s_1) + \log p(a_1 | s_1, \theta) \quad (2b)$$

$$+ \sum_{t=2}^T [\log p(s_t | s_{t-1}, a_{t-1}) + \log p(a_t | s_t, \theta)] \quad (2c)$$

$$= \sum_{t=1}^T [\log p(a_t | s_t, \theta)] + A \quad (2d)$$

216 A are all the terms that do not contain θ . Substituting $2d$ in $1f$:

$$\begin{aligned} 217 \log p(\tau) &\geq \mathbb{E}_{q(z|\tau)} \left[\mathbb{E}_{p(\theta|z)} [\log(p(\tau | \theta))] \right] - \text{KL}(q(z | \tau) \| p(z)) \\ 218 &= \mathbb{E}_{q(z|\tau)} \left[\mathbb{E}_{p(\theta|z)} \left[\sum_{t=1}^T [\log p(a_t | s_t, \theta)] \right] \right] + A - \text{KL}(q(z | \tau) \| p(z)) \end{aligned} \quad (3)$$

222 We ignore the terms in A as these cannot be subject to maximization as we do not have access to
223 the state transition probabilities (although this can be modeled with a world model, we leave it as a
224 direction for future work).

225 Therefore, our modified ELBO is:

$$226 \quad mELBO = \sum_{t=1}^T \left[\mathbb{E}_{q(z|\tau)} \left[\mathbb{E}_{p(\theta|z)} [\log p(a_t | s_t, \theta)] \right] \right] - \text{KL}(q(z | \tau) \| p(z)) \quad (4)$$

230 3.2 VARIATIONAL AUTOENCODER FOR POLICIES

232 Since we now have a modified ELBO objective, we shall now try to approximate its components
233 with a variational autoencoder. Let ϕ_{enc} be the parameters of the VAE encoder that variationally
234 maps trajectories to z , and ϕ_{dec} be the parameters of the VAE decoder. We assume the latent z is
235 distributed with mean zero and unit variance. We construct the VAE decoder to approximate $p(\theta | z)$
236 with $p_{\phi_{dec}}(\theta | z)$. Since $a_t \sim \mathcal{N}(\pi(s_t, \theta), \sigma^2)$:

$$237 \quad p(a_t | s_t, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(a_t - \pi(s_t, \theta))^2}{2\sigma^2}\right) \quad (5)$$

240 Our objective is to maximize the $mELBO$. The likelihood of trajectory $\tau_k = \{s_t^k, a_t^k\}_{t=1}^T$ for the
241 given VAE parameters is:

$$\begin{aligned} 243 \mathcal{L}(\{s_t^k, a_t^k\}_{t=1}^T | \phi_{enc}, \phi_{dec}) &= \sum_{t=1}^T \mathbb{E}_{q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T)} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta | z)} [\log p(a_t^k | s_t^k, \theta)] \right] \\ 244 &\quad - \text{KL}(q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T) \| p(z)) \\ 245 &= C - \frac{1}{2\sigma^2} \sum_{t=1}^T \mathbb{E}_{q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T)} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta | z)} [(a_t^k - \pi(s_t^k, \theta))^2] \right] \\ 246 &\quad - \text{KL}(q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T) \| p(z)) \end{aligned} \quad (6)$$

252 For computational stability, we construct our decoder to be a deterministic function $f_{\phi_{dec}}$, i.e.,
253 $p_{\phi_{dec}}(\theta | z)$ becomes $\delta(\theta - f_{\phi_{dec}}(z))$, therefore:

$$\begin{aligned} 254 \mathcal{L}(\{s_t^k, a_t^k\}_{t=1}^T | \phi_{enc}, \phi_{dec}) &= C - \frac{1}{2\sigma^2} \sum_{t=1}^T \mathbb{E}_{q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T)} [(a_t^k - \pi(s_t^k, f_{\phi_{dec}}(z)))^2] \\ 255 &\quad - \text{KL}(q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T) \| p(z)) \end{aligned}$$

259 Enforcing $p(z) = \mathcal{N}(0, I)$, and ignoring constants unaffected by the VAE parameters, we get:

$$\begin{aligned} 261 \mathcal{L}(\{s_t^k, a_t^k\}_{t=1}^T | \phi_{enc}, \phi_{dec}) &= - \sum_{t=1}^T \mathbb{E}_{q_{\phi_{enc}}(z | \{s_t^k, a_t^k\}_{t=1}^T)} [(a_t^k - \pi(s_t^k, f_{\phi_{dec}}(z)))^2] \\ 262 &\quad - \beta_{kl} \sum_{i=1}^{\dim(z)} (\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2) \end{aligned} \quad (7)$$

267 where, $(\mu, \sigma) = f_{\phi_{enc}}(\{s_t^k, a_t^k\}_{t=1}^T)$, $z \sim \mathcal{N}(\mu, \sigma)$, and β_{kl} is the regularization weight. Since
268 the decoder in the VAE outputs the parameter of a secondary network, we shall use a conditional
269 Hypernetwork, specifically the one that was developed for continual learning by von Oswald et al.
(2020).

270 3.3 POLICY DIFFUSION

271
272 In practice, we see that approximating $p(z) = \mathcal{N}(0, I)$ is suboptimal, and therefore we set β_{kl} to
273 a very small number $\sim (10^{-9}, 10^{-6})$. After training the VAE to maximize the objective provided
274 in Equation 7 with this β_{kl} , we have access to this latent space z and can train a diffusion model to
275 learn its distribution $p(z)$. The diffusion process in the latent space involves gradually adding noise
276 to the latent variable $\mathbf{z}_0 = \mathcal{E}(\tau)$ over a sequence of time steps $t \in \{0, 1, \dots, T\}$. This process can
277 be described by a forward noising process:

$$278 \epsilon_t = q(z_t | z_{t-1}) = \mathcal{N}(z_t; z_{t-1}\sqrt{1 - \beta_t}, \beta_t \mathbf{I}) \quad (8)$$

280 making the forward process Markovian. The reverse or generative process $p_{\phi_{dif}}(z_T)$ reverts noise
281 from an isotropic Gaussian into a sample z_0 in $q(\mathbf{z})$ and contains a similar Markov structure.

$$282 \quad 283 \quad 284 \quad 285 \quad 286 \quad 287 \quad 288$$

$$p_{\phi_{dif}}(z_0) = p(z_T) \prod_{t=1}^T p_{\phi_{dif}}(z_{t-1} | z_t), \quad p_{\phi_{dif}}(z_{t-1} | z_t) = \mathcal{N}(z_{t-1}; \mu_{\phi_{dif}}(z_t, t), \Sigma_{\phi_{dif}}(z_t, t)) \quad (9)$$

289 We can condition the latent denoising process on the current state and/or the task identifier c of the
290 policy required. Therefore the model shall be approximating $p_{\phi_{dif}}(z_{t-1} | z_t, c)$. After denoising for
291 a given state and task identifier, we can convert the denoised latent to the required policy.

292 Therefore, to sample from $p(\theta)$, first sample z using the trained diffusion model $z \sim p_{\phi_{dif}}(z_0)$, and
293 then apply the deterministic function $f_{\phi_{dec}}$ to the sampled z .

294 4 EXPERIMENTS

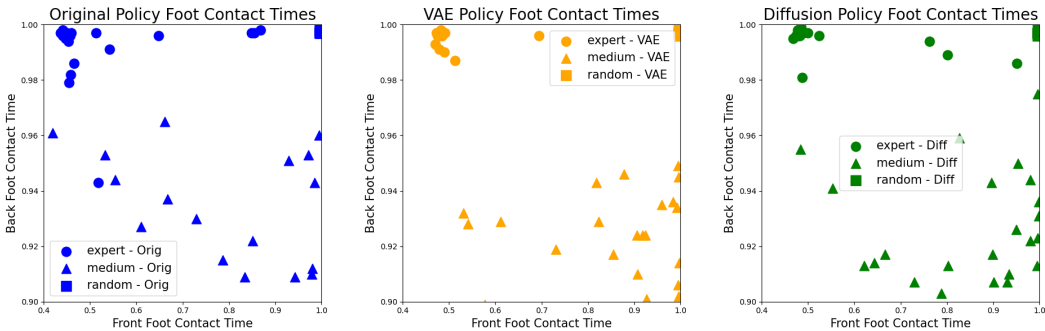
295
296 We first analyze the behavior reconstruction of LWD, followed by the effect of breaking down
297 the trajectory into shorter snippets. Then, we perform an ablation over the size of LWD, showing
298 that a larger model can mitigate the problems introduced by snipping trajectories. After this, we
299 benchmark LWD on the MT10 suite of tasks in Metaworld, showcasing multitask advantages, as
300 well as a benchmark on human-generated data in the PushT environment, showcasing its ability to
301 run at longer action horizons. All our numerical results are shown across three seeds.

302 We focus on demonstrating results in state-based low-dimension observation spaces. Thus, the
303 generated policies are always Multi-Layer Perceptrons (MLP) with 2 hidden layers with 256 neurons
304 each. In the VAE, the encoder is a sequential network that flattens the trajectory and compresses it
305 to a low-dimension latent space. The decoder is a conditional hypernetwork from the hypernettorch
306 package (Ehret et al., 2021). For the diffusion model, we adapt the model provided by Rombach
307 et al. (2022a). For all experiments the latent space is \mathbb{R}^{256} , the KL regularization weight $\beta_{kl} = 10^{-8}$,
308 the learning rate is 10^{-4} with the Adam optimizer.

309 4.1 BEHAVIOR RECONSTRUCTION ANALYSIS

310
311 Here, we ask – how does LWD perform in reconstructing the behavior of the original policies that
312 generated the trajectory dataset? Is it able to reproduce different behaviors for the same task?

313
314 First, we analyze the behavior reconstruction capability of different components of LWD. For this
315 experiment, we use the D4RL (Fu et al., 2020) halfcheetah dataset. Each trajectory in this dataset
316 has a length of 1000. We combine trajectory data from three original behavior policies provided in
317 this dataset: expert, medium, and random. Following Batra et al. (2023), we track the foot contact
318 timings of each trajectory as a metric for measuring behavior. For each behavior policy, we get 32
319 trajectories. These timings are normalized to the trajectory length and are shown in Figure 3. For
320 each plot, the x-axis denotes the foot contact percentage of the front foot, while the y-axis denotes
321 the foot contact percentage of the back foot. We first visualize the foot contact timings of the original
322 policies in Figure 3a. Then, we train the VAE model on this dataset to embed our trajectories into
323 a latent space. We then apply the hypernetwork decoder to generate policies from these latents.
These policies are then executed on the halfcheetah environment, to create trajectories. We plot the



(a) Original policies that provide the trajectory dataset. (b) VAE generated policies from trajectories. (c) VAE + diffusion generated policies.

Figure 3: **Foot-contact times shown for various trajectories on the Half Cheetah task.** Figure 3a: We use foot contact times as the chosen metric to show different behaviors for the half cheetah run task by different policies. Figure 3b: Our VAE can embed these behaviors into a latent space and then reconstruct a policy from them with the same behavioral patterns. Figure 3c Our diffusion model then learns the distribution of this latent space. We train it with the task label as the conditioning.

foot contact timings of these generated policies in Figure 3b. We see that the VAE captures each of the original policy’s foot contact distributions, therefore empirically showing that the assumption $p_{\phi_{dec}}(\theta | z) = \delta(\theta - f_{\phi_{dec}}(z))$ is reasonable. Then, we train a latent diffusion model conditioned on a behavior specifier (i.e., one task ID per behavior). In Figure 3c, we show the distribution of foot contact percentages of the policies generated by the behavior specifier conditioned diffusion model. We see that the diffusion model can learn the conditional latent distribution well, and the behavior distribution of the decoded policies of the sampled latent matches the original distribution. Note that to sample policies during inference, we do not need to encode trajectories; rather, we need to sample latents using the diffusion model and use the hypernetwork decoder to decode a policy from it.

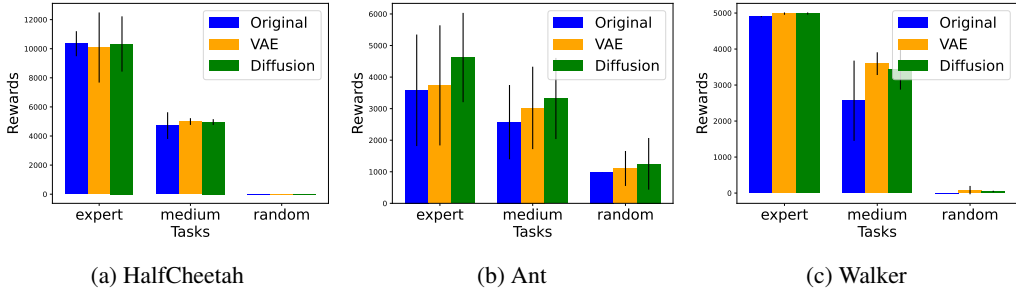


Figure 4: **Reconstruction Rewards:** For each of the 3 environments shown above, the generated policy from trajectory decoded VAE and task-conditioned diffusion model, achieves similar total objective as the original policies. Each bar indicates the mean total objective obtained with error lines denoting the standard deviation.

Another way to analyze the behavior reconstruction capability of LWD is to compare the rewards obtained during a rollout. Figure 4 shows us the total objective obtained by the original, VAE-decoded, and diffusion-denoised policies. We see that the VAE-decoded and diffusion-generated policies achieve similar rewards to the original policy for each behavior.

Apart from these plots, we use Jensen-Shannon divergence to quantify the difference between two distributions of foot contact timings. Table 1 shows the JS divergence between the empirical distribution of the foot contact timings of the original policies and those generated by LWD. The lower this value is, the better. As a metric to capture the stochasticity in the policy and environment, we get the JS divergence between two successive sets of trajectories generated by the same original policy, which we shall denote SOS (Same as source). A policy having a JS divergence score lesser than this value indicates that that policy is indistinguishable from the original policy by behavior. As a

378 baseline for this experiment, we train a large (5-layer, 512 neurons each) behavior-conditioned MLP
 379 on the same mixed dataset with MSE loss. We see that policies generated by LWD consistently
 380 achieve a lower JS divergence score than the MLP baseline for expert and medium behaviors. The
 381 random behavior is difficult to capture as the actions are almost Gaussian noise. Surprisingly, for
 382 the HalfCheetah environment, policies generated by LWD for expert and medium had lower scores
 383 than SOS, making it behaviorally indistinguishable from the original policy.
 384

Environment	Source Policy	Target Policy		
		SOS	MLP	LWD
Ant	Expert	0.187 ± 0.142	1.272 ± 0.911	0.510 ± 0.159
	Medium	0.624 ± 0.232	1.907 ± 0.202	1.328 ± 0.283
	Random	1.277 ± 1.708	4.790 ± 0.964	8.859 ± 0.792
HalfCheetah	Expert	0.158 ± 0.146	2.810 ± 1.139	0.088 ± 0.050
	Medium	0.275 ± 0.196	0.692 ± 0.787	0.194 ± 0.157
	Random	0.0467 ± 0.009	0.11 ± 0.009	0.104 ± 0.0187
Walker2D	Expert	0.342 ± 0.329	2.879 ± 1.493	1.093 ± 0.310
	Medium	0.078 ± 0.058	0.165 ± 0.126	0.155 ± 0.091
	Random	0.080 ± 0.004	60.514 ± 52.461	2.776 ± 1.260

385
386
387
388
389
390
391
392
393
394
395
396 Table 1: **Behavior Reconstruction:** JS divergence between foot contact distributions from source
 397 and target policies. The lower the value, the better.

398
399 4.2 ENCODING TRAJECTORY SNIPPETS

400
401 Here, we ask the question – can LWD generate policies that are faithful to the original policies,
 402 even when provided with only a snippet of the trajectory data? For most robotics use cases, it
 403 is impossible to train on long trajectories due to the computational limitations of working with
 404 large batches of long trajectories. Therefore, we analyze the effect of sampling smaller sections
 405 of trajectories from the dataset. After training a VAE for the D4RL halfcheetah dataset on three
 406 policies (expert, medium, and random), we encode all the trajectories in the mixed dataset to the
 407 latent space. We then perform Principal Component Analysis (PCA) on this set of latents and select
 408 the first two principal components. Figure 5a shows us a visualization of this latent space. We see
 409 that the VAE has learned to encode the three sets of trajectories to be well separable. Next, we run
 410 the same experiment, but now we sample trajectory snippets of length 100 from the dataset instead
 411 of the full-length (1000) trajectories. Figure 5b shows us the PCA on the encoded latents of these
 412 trajectory snippets. We see that the separability is now harder in the latent space. Surprisingly, we
 413 noticed that after training our VAE on the snippets, the decoded policies from randomly snipped
 414 trajectories were still faithfully behaving like their original policies. We believe that this is because
 415 of the cyclic nature of the halfcheetah task, and all trajectory snippets have enough information to
 416 indicate its source policy.

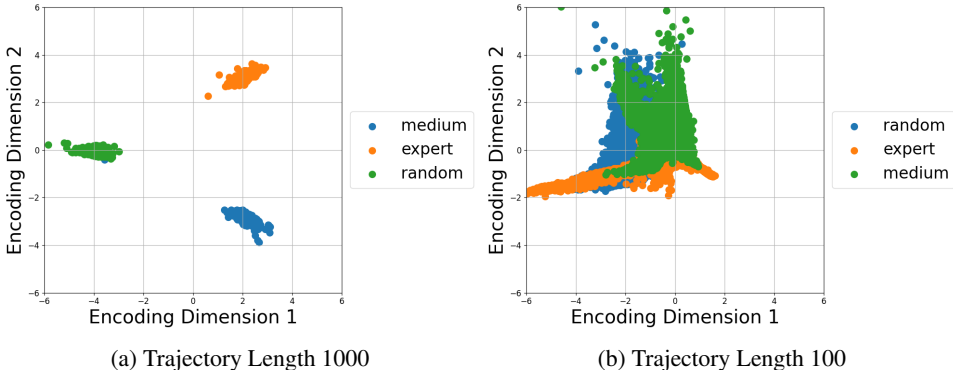
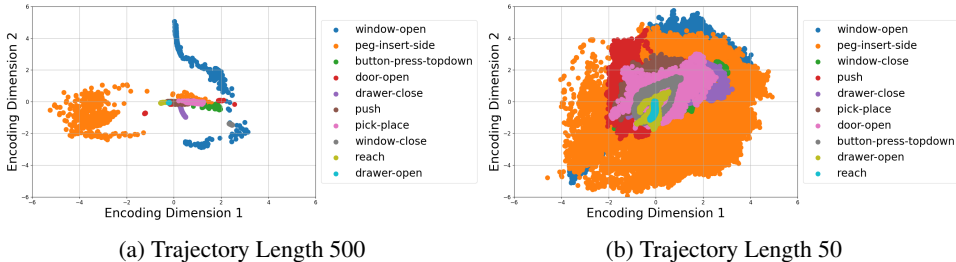


Figure 5: **Effect of trajectory snipping** in HalfCheetah. Top two principal components of the latent.

To validate this hypothesis, we analyze our method on trajectory snippets for non-cyclic tasks. We choose the MT10 suite of tasks in Metaworld (Yu et al., 2020). We utilize the hand-crafted expert

432 policy for each of the tasks in MT10 to collect trajectory data. For each task, we collect 1000
 433 trajectories of length 500.
 434



444 Figure 6: **Effect of trajectory snipping** in MT10. Top two principal components of the latent.
 445

446 Figure 6a shows the principal components of the latents of the full trajectories in the dataset,
 447 and Figure 6b shows the same for the split trajectories. We can see that the separability of
 448 different tasks is much harder in this case. Further, we noticed that the decoded policies from
 449 the trajectory snippets did not perform as well as the original policies - for the same decoder
 450 size as the half cheetah task. However, we did notice an ability to generate policies that
 451 were faithful to the original policies when we increased the decoder size. This is discussed
 452 next, in subsection 4.3.
 453
 454
 455
 456
 457
 458

Policy	MLP512	MLP256	MLP128	LWD
# Parameters	1396.2k	370.4k	103.3k	77.1k
button press	1.00	1.00	1.00	1.00
door open	1.00	1.00	1.00	1.00
drawer close	1.00	0.87	0.87	1.00
drawer open	1.00	1.00	0.67	1.00
peg insert	0.13	0.20	0.27	0.53
pick place	0.13	0.00	0.00	0.13
push	0.13	0.00	0.13	0.47
reach	0.53	0.33	0.267	0.47
window close	1.00	1.00	1.00	1.00
window open	1.00	1.00	1.00	1.00
Mean over seeds	0.693	0.640	0.620	0.76
Stddev over seeds	0.072	0.04	0.061	0.052

459 4.3 VAE DECODER SIZE ABLATION
 460

461 As noted in subsection 4.2, the size of the hypernetwork decoder influences the quality of
 462 decoded policies for the MT10 task suite, when trained on trajectory snippets. Here we conduct
 463 an ablation on trajectory snippets. Here we conduct
 464 an ablation on the decoder size, evaluating the average success rate of decoded policies across all
 465 MT10 tasks. Figure 7 illustrates the performance of decoders with varying sizes, denoted as xs
 466 ($3.9M$ parameters), s ($7.8M$ parameters), m ($15.6M$ parameters), and l ($31.2M$ parameters). It’s
 467 important to note that despite the substantial parameter count of the hypernetwork decoder, the re-
 468 sulting inferred policy remains relatively small ($< 100K$ parameters, see Table 2). The results
 469 demonstrate that increasing the decoder size consistently improves the average success rate of the
 470 decoded policies.
 471

472 This contrasts with the observations from the HalfCheetah environment, where even smaller de-
 473 coders generated accurate policies from trajectory snippets. We hypothesize that this discrepancy
 474 stems from two key factors. First, the cyclic nature of HalfCheetah provides sufficient information
 475 within the snippets to infer the underlying policy. Second, the increased complexity of the MT10
 476 tasks means that snippets may lack crucial information for policy inference. For instance, in a
 477 pick-and-place task, a snippet might only capture the “pick” action, leaving the latent representation
 478 without sufficient information to infer the “place” action.
 479

480 4.4 METAWORLD MT10 SIZE BENCHMARK

481 In the previous experiments, we have shown that LWD can learn a distribution of policies from a
 482 trajectory dataset. However, another common use case for robotic learning is multi-task imitation
 483 learning. In this experiment, we study the ability of LWD to learn a task conditional distribution of
 484 policies. We test the performance of LWD on the Metaworld MT10 benchmark and compare it to
 485 a vanilla multi-task MLP policy trained on the same dataset. The baselines are three MLP policies,
 each with 5 hidden layers of the same size (512, 256 and 128 perceptions). As seen in Table 2,
 LWD outperforms the vanilla multi-task policy in the average success rate across all tasks. We also

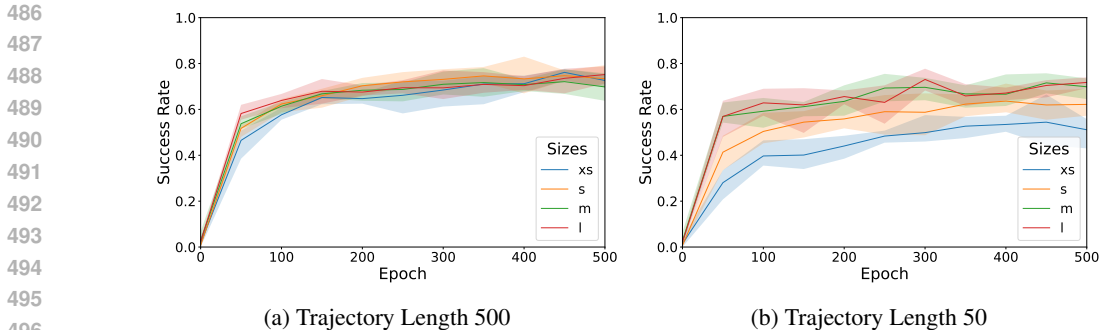


Figure 7: **Effect of VAE decoder size:** For longer trajectories, even the smallest decoder, *xs*, is sufficient to give high task performance. For shorter trajectories, a larger decoder model helps maintain the same level of performance.

see that the policy generated by LWD has a smaller parameter count than the vanilla multi-task policy, with about 1/18th the number of parameters.

4.5 EFFECT OF ACTION HORIZON

The previous experiments have shown that LWD can learn a distribution of behaviors for a task or a distribution of policies for a set of tasks. All these experiments were conducted with a single policy being generated at the start of the episode. However, for environments where the training dataset has a high variance, it is important to diffuse out locally optimal policies. That is, we need to diffuse a policy out every H_a steps, where H_a is the action horizon. This is what we observe for Diffusion Policy (Chi et al., 2024) as well. We use their PushT task, whose dataset has fewer trajectories with a higher variance. We compare the performance of LWD to Diffusion Policy on this task. Specifically, we compare the performance change as we change the action horizon from 8 to 128, relative to the best value. We now condition our latent diffusion model on the current state. We see in Figure 8 that Diffusion Policy had the best performance at a horizon of 16, whereas LWD had the best performance at a horizon of 8. But note that as the action horizon increases, the performance of Diffusion Policy degrades much faster than LWD. Therefore, for a relative performance reduction of $\sim 25\%$, LWD requires 1/4th the number of expensive diffusion model queries.

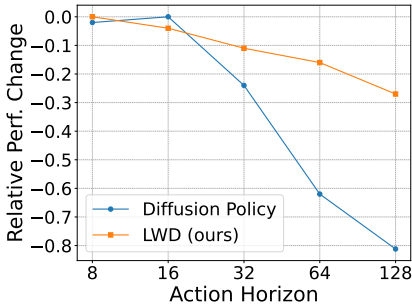


Figure 8: LWD maintains performance for longer action horizons, allowing fewer diffusion queries for the same episode length.

5 LIMITATIONS AND FUTURE WORK

Although LWD has promising results, it has some limitations and avenues for future work. LWD needs a lot of trajectory data to learn the distribution of policies accurately. Moreover, by virtue of generating closed-loop policies, LWD is more prone to see out-of-distribution states when compared to methods that diffuse multi-step trajectories. These limitations serve as potential directions for future work. For example, it would be interesting to investigate methods to learn the distribution of policies from fewer demonstration trajectories. Another avenue is to use our method with a hypernetwork that can output weights for a transformer network, to take sequence data as input, or a vision transformer network, to take image data as input.

6 CONCLUSION

We present LWD, a method to learn a distribution of policies from a heterogeneous set of demonstration trajectories. We first embed trajectories into a latent space and then learn the distribution of policies in this latent space. We then decode these latents to generate policies using a hypernetwork decoder. We show that LWD can reproduce the original policies present in the demonstration trajectories, in two cases. First, we show that we can reproduce multiple behaviors for the same task. We also show that we can reproduce policies used for multi-task learning. Finally, we discuss how

540 LWD can be used with high-variance data, and compare it to baselines. We believe that LWD can be
541 a useful tool for learning from demonstration, and can be used in a variety of applications, including
542 robotics, reinforcement learning, and imitation learning.
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

- 594
595
596 Sumeet Batra, Bryon Tjanaka, Matthew C Fontaine, Aleksei Petrenko, Stefanos Nikolaidis, and
597 Gaurav Sukhatme. Proximal policy gradient arborescence for quality diversity reinforcement
598 learning. *arXiv preprint arXiv:2305.13795*, 2023.
- 599 Luca Bertinetto, João F. Henriques, Jack Valmadre, Philip Torr, and Andrea Vedaldi. Learn-
600 ing feed-forward one-shot learners. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and
601 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran
602 Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper_files/
603 paper/2016/file/839ab46820b524afda05122893c2fe8e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/839ab46820b524afda05122893c2fe8e-Paper.pdf).
- 604 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn,
605 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian
606 Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalash-
607 nikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deek-
608 sha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez,
609 Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi,
610 Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent
611 Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and
612 Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In *arXiv preprint
613 arXiv:2212.06817*, 2022.
- 614 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choro-
615 manski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu,
616 Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alex Her-
617 zog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov,
618 Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Hen-
619 ryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo,
620 Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut,
621 Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart,
622 Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-
623 2: Vision-language-action models transfer web knowledge to robotic control. In *arXiv preprint
624 arXiv:2307.15818*, 2023.
- 625 Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake,
626 and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The Inter-
627 national Journal of Robotics Research*, 2024.
- 628 Open X-Embodiment Collaboration, Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhiram Mad-
629 dukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay
630 Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khaz-
631 atsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg,
632 Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh
633 Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim,
634 Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea
635 Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher
636 Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne
637 Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov,
638 Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao,
639 Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan,
640 Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao
641 Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki
642 Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Ra-
643 dosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters,
644 Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen
645 Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon
646 Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan
647 Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan
Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken
Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin

- 648 Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan
649 Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yun-
650 liang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca
651 Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina,
652 Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong
653 Zhang, Mingyu Ding, Minh Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki
654 Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Nor-
655 man Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani,
656 Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David
657 Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundareshan, Qiuyu Chen, Quan
658 Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Bajjal, Rosario
659 Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah,
660 Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry
661 Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun
662 Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany,
663 Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel
664 Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya
665 Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev,
666 Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vin-
667 cent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong
668 Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao
669 Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying
670 Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen
671 Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang
672 Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen
673 Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X
674 models. <https://arxiv.org/abs/2310.08864>, 2023.
- 674 Coline Devin, Abhishek Gupta, Trevor Darrell, P. Abbeel, and Sergey Levine. Learning mod-
675 ular neural network policies for multi-task and multi-robot transfer. *2017 IEEE International
676 Conference on Robotics and Automation (ICRA)*, pp. 2169–2176, 2016. URL [https://api.
677 semanticscholar.org/CorpusID:18015872](https://api.semanticscholar.org/CorpusID:18015872).
- 678 Benjamin Ehret, Christian Henning, Maria R. Cervera, Alexander Meulemans, Johannes von Os-
679 wald, and Benjamin F. Grewe. Continual learning in recurrent neural networks. In *International
680 Conference on Learning Representations*, 2021. URL [https://arxiv.org/abs/2006.
681 12109](https://arxiv.org/abs/2006.12109).
- 682 Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian
683 Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In Aleks-
684 sandra Faust, David Hsu, and Gerhard Neumann (eds.), *Proceedings of the 5th Conference on
685 Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pp. 158–168. PMLR,
686 08–11 Nov 2022. URL [https://proceedings.mlr.press/v164/florence22a.
687 html](https://proceedings.mlr.press/v164/florence22a.html).
- 688 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
689 data-driven reinforcement learning, 2020.
- 690 David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- 691 Shashank Hegde, Sumeet Batra, KR Zentner, and Gaurav Sukhatme. Generating behaviorally di-
692 verse policies with latent diffusion models. *Advances in Neural Information Processing Systems*,
693 36:7541–7554, 2023.
- 694 Shashank Hegde, Zhehui Huang, and Gaurav S Sukhatme. Hyperppo: A scalable method for find-
695 ing small policies for robotic control. In *2024 IEEE International Conference on Robotics and
696 Automation (ICRA)*, pp. 10821–10828. IEEE, 2024.
- 697 Negin Heravi, Ayzaan Wahid, Corey Lynch, Peter R. Florence, Travis Armstrong, Jonathan Tomp-
698 son, Pierre Sermanet, Jeannette Bohg, and Debidatta Dwibedi. Visuomotor control in multi-object
699 scenes using object-aware representations. *2023 IEEE International Conference on Robotics and
700 Automation (ICRA)*, pp. 10821–10828. IEEE, 2024.
- 701

- 702 *Automation (ICRA)*, pp. 9515–9522, 2022. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:248798881)
703 [CorpusID:248798881](https://api.semanticscholar.org/CorpusID:248798881).
704
- 705 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In
706 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-
707 Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Con-*
708 *ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,*
709 *2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html)
710 [4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html).
- 711 Xiaoyu Huang, Yufeng Chi, Ruofeng Wang, Zhongyu Li, Xue Bin Peng, Sophia Shao, Borivoje
712 Nikolic, and Koushil Sreenath. Diffuseloco: Real-time legged locomotion control with diffusion
713 from offline datasets. *arXiv preprint arXiv:2404.19264*, 2024.
- 714 Stephen James, Z. Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot learning
715 benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5:3019–3026, 2019.
716 URL <https://api.semanticscholar.org/CorpusID:202889132>.
717
- 718 Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter net-
719 works. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Ad-*
720 *vances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.,
721 2016. URL [https://proceedings.neurips.cc/paper_files/paper/2016/](https://proceedings.neurips.cc/paper_files/paper/2016/file/8bf1211fd4b7b94528899de0a43b9fb3-Paper.pdf)
722 [file/8bf1211fd4b7b94528899de0a43b9fb3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/8bf1211fd4b7b94528899de0a43b9fb3-Paper.pdf).
- 723 Justin Lovelace, Varsha Kishore, Chao Wan, Eliot Shekhtman, and Kilian Q Weinberger. Latent
724 diffusion for language generation. *Advances in Neural Information Processing Systems*, 36, 2024.
725
- 726 Utkarsh Aashu Mishra, Shangjie Xue, Yongxin Chen, and Danfei Xu. Generative skill chaining:
727 Long-horizon skill planning with diffusion models. In *Conference on Robot Learning*, 2023.
728 URL <https://api.semanticscholar.org/CorpusID:261685884>.
- 729 Xue Peng, Erwin Coumans, Tingnan Zhang, Tsang-Wei Lee, Jie Tan, and Sergey Levine. Learning
730 agile robotic locomotion skills by imitating animals. 07 2020. doi: 10.15607/RSS.2020.XVI.064.
731
- 732 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
733 resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer*
734 *Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp.
735 10674–10685. IEEE, 2022a. doi: 10.1109/CVPR52688.2022.01042. URL [https://doi.](https://doi.org/10.1109/CVPR52688.2022.01042)
736 [org/10.1109/CVPR52688.2022.01042](https://doi.org/10.1109/CVPR52688.2022.01042).
- 737 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
738 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
739 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022b.
- 740 Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for
741 robotic manipulation. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.
742
- 743 Wenhui Tan, Bei Liu, Junbo Zhang, Ruihua Song, and Jianlong Fu. Multi-task manipulation policy
744 modeling with visuomotor latent diffusion. *ArXiv*, abs/2403.07312, 2024. URL [https://api.](https://api.semanticscholar.org/CorpusID:268364239)
745 [semanticscholar.org/CorpusID:268364239](https://api.semanticscholar.org/CorpusID:268364239).
- 746 Julen Urain, Niklas Funk, Jan Peters, and Georgia Chalvatzaki. Se(3)-diffusionfields: Learning
747 smooth cost functions for joint grasp and motion optimization through diffusion. *2023 IEEE*
748 *International Conference on Robotics and Automation (ICRA)*, pp. 5923–5930, 2022. URL
749 <https://api.semanticscholar.org/CorpusID:252367206>.
- 750 Johannes von Oswald, Christian Henning, Benjamin F. Grewe, and João Sacramento. Continual
751 learning with hypernetworks. In *International Conference on Learning Representations*, 2020.
752 URL <https://arxiv.org/abs/1906.00695>.
753
- 754 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
755 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.

756 Yanjie Ze, Ge Yan, Yueh-Hua Wu, Annabella Macaluso, Yuying Ge, Jianglong Ye, Nicklas Hansen,
757 Li Erran Li, and Xiaolong Wang. Multi-task real robot learning with generalizable neural feature
758 fields. *CoRL*, 2023.
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809