Latent Weight Diffusion: Generating reactive policies instead of trajectories

Shashank Hegde

University of Southern California khegde@usc.edu

Satyajeet Das

University of Southern California

Gautam Salhotra (Google) Intrinsic LLC

Gaurav S. Sukhatme University of Southern California

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. This is especially an issue for robotic tasks that require high control frequency. Further, there is a known trade-off between performance and action horizon for Diffusion Policy (DP), a popular model for generating trajectories: fewer diffusion queries accumulate greater trajectory tracking errors. For these reasons, 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 and a world model to generate closed-loop policies (weights for neural policies) for robotic tasks, rather than generating trajectories. Learning the behavior distribution through parameter space over trajectory space offers two key advantages: longer action horizons (fewer diffusion queries) & robustness to perturbations while retaining high performance; and a lower inference compute cost. To this end, we show that LWD has higher success rates than DP when the action horizon is longer and when stochastic perturbations exist in the environment. Furthermore, LWD achieves multitask performance comparable to DP while requiring just $\sim 1/45$ th of the inference-time FLOPS per step.

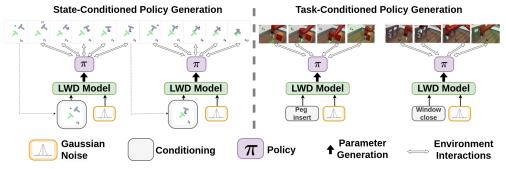


Figure 1: Latent Weight Diffusion (LWD) generates policies from heterogeneous trajectory data. With state-conditioned policy generation, the diffusion model can run inference at a lower frequency. With task-conditioned policy generation, the generated policies can be small yet maintain task-specific performance. Demonstrations of this work can be found on the project website: https://sites.google.com/view/lwd2025/home.

39th Conference on Neural Information Processing Systems (NeurIPS 2025) Workshop: Embodied World Models for Decision Making.

1 Introduction

The rise of open-source robotic datasets has made imitation learning a promising approach for robotic manipulation and locomotion tasks [10, 43]. While methods like Behavioral Cloning [13] and transformer-based models (e.g., RT-1 [5]) have shown promise, they struggle with multimodal action distributions. For example, in navigation tasks where both "turn left" and "turn right" are valid, these models often predict an averaged action, i.e., "go straight", leading to suboptimal performance.

Diffusion models offer a compelling alternative, providing continuous outputs and learning multimodal action distributions [51]. Action trajectory diffusion for robotic tasks [9] has shown promise but incurs high computational costs, particularly at high control frequencies. Moreover, such trajectory diffusion models are susceptible to the trade-off between performance and action horizon (or action chunk size, representing the number of environment interactions between consecutive trajectory generations). Fewer diffusion queries lead to larger action chunks, giving greater trajectory tracking errors.

To overcome these limitations, we introduce Latent Weight Diffusion (LWD), a novel approach that uses latent diffusion and a world model to generate closed-loop policies directly in parameter space, bypassing trajectory generation. LWD first encodes demonstration trajectories into a latent space, then learns their distribution using a diffusion model, and finally decodes them into policy weights via a hypernetwork [15]. The generated policy is also optimized with model-based imitation learning using a co-trained world (dynamics) model [16], which helps in understanding the environment transitions during training. This approach leverages the success of latent diffusion techniques in vision [47] and language [34], and combines them with learned dynamics models, bringing their advantages to robotic control. The world model, and accompanying loss terms, help the agent learn the optimal policy that can be backpropagated through the learned (differentiable) dynamics, and also apply corrective actions to bring back the agent states back into distribution of the input trajectory dataset. For LWD, the action horizon corresponds to the number of environment interactions between consecutive policy weight generations. To achieve trajectory encoding and policy parameter decoding, we derive a novel objective function described in Section 3.1, and show that we can approximate its components with a hypernetwork-based VAE and a World Model, and optimize it using a novel loss function described in Section 3.2. This paper provides the following key contributions:

- 1. **Theoretical Foundations for generating policies**: We derive a novel objective function, which when optimized, allows us to generate policy parameters instead of action trajectories. For this, we integrate concepts from latent diffusion, hypernetworks, and world models.
- 2. Longer Action Horizons & Robustness to Perturbations: By generating closed-loop policies under learned dynamics, LWD mitigates trajectory tracking errors, enabling policies to operate over extended time horizons with fewer diffusion queries. Additionally, closed-loop policies are reactive to environmental changes, ensuring LWD-generated policies remain robust under stochastic disturbances.
- 3. **Lower Inference Costs**: The computational burden of generalization is shifted to the diffusion model, allowing the generated policies to be smaller and more efficient.

We validate these contributions through experiments on the PushT task [9], the Lift and Can tasks from Robomimic [36], and 10 tasks from Metaworld [58]. On Metaworld, achieves comparable performance to Diffusion Policy but with a $\sim 45 x$ reduction in FLOPS per step, representing a significant improvement in computational efficiency (FLOPS per step are the floating point operations per second, amortized over all steps of the episode). Analysis across a range of benchmark robotic locomotion and manipulation tasks, demonstrates LWD's ability to accurately capture the *behavior distribution* of diverse trajectories, showcasing its capacity to learn a distribution of behaviors.

2 Related Work

2.1 Imitation Learning and Diffusion for Robotics

Behavioral cloning has advanced with transformer-based models like PerAct [49] and RT-1 [5], achieving strong task performance. Vision-language models (e.g., RT-2 [4]) extend this by interpreting text as actions, while RT-X [10] generalizes across robot embodiments. Object-aware representations [24],

energy-based models, and temporal abstraction (e.g., implicit behavioral cloning [13], sequence compression [61]) further improve multitask learning. DBC [7] enhances robustness to sensor noise, whereas LWD targets environmental perturbations affecting system dynamics, such as object shifts or execution-time disturbances.

Alongside these advances, diffusion models (originally introduced for generative modeling [25, 48]) have emerged as powerful tools for robotics. Trajectory-based diffusion approaches capture multimodal action distributions [9], while goal-conditioned methods such as BESO [46] and Latent Diffusion Planning [30] improve efficiency through latent-space conditioning. Diffusion has also been leveraged for grasping and motion planning [52, 35, 6], skill chaining [39], and locomotion control [28]. Hierarchical extensions such as ChainedDiffuser [57], SkillDiffuser [33], and multitask latent diffusion [51] address long-horizon action planning, though trajectory tracking remains challenging. Recently, OCTO [40] demonstrates that diffusion-based generalist robot policies.

2.2 Hypernetworks and Policy Generation

Hypernetworks, first introduced by [15], are neural architectures designed to estimate parameters for secondary networks, finding applications across various domains. Initially applied to metalearning scenarios for one-shot learning tasks [3], hypernetworks have recently been extended to robot policy representations [23]. Their approach aligns conceptually with Dynamic Filter Networks [29], emphasizing dynamic adaptability to input data.

Recent developments also integrate Latent Diffusion Models (LDMs) into hypernetwork-like settings, modeling training dynamics in parameter spaces [41]. Specifically, LDMs have been successfully used to generate behavior-conditioned policies from textual descriptions [22] and trajectory embeddings [32], as well as to learn distributions over complex model architectures like ResNet [54]. Notably, while these recent approaches [22, 32] depend on pre-collected policy datasets, the current work distinguishes itself by removing this requirement.

2.3 World Models

Ha and Schmidhuber [16] introduced world models for forecasting in latent space. PlaNet [18] extended this with pixel-based dynamics learning and online planning. Dreamer [17] learned latent world models with actor-critic RL for long-horizon behaviors, followed by DreamerV2 [19] with discrete representations achieving human-level Atari, and DreamerV3 [20] scaling robustly across domains. IRIS [37] used transformers for sequence modeling, reaching superhuman Atari in two hours. SLAC [31] showed that stochastic latent variables accelerate RL from high-dimensional inputs. VINs [50] embedded differentiable value iteration for explicit planning, while E2C [55] combined VAEs with locally linear dynamics. DayDreamer [56] enabled real robot learning in one hour, and MILE [27] adapted Dreamer to CARLA with 31% gains. Model-based imitation learning was further scaled to large self-driving datasets by Popov et al. [44]. Recent advances include SafeDreamer [60] for safety, STORM [38] with efficient transformers, UniZero [59] for joint optimization of models and policies, and Time-Aware World Models [8] incorporating temporal dynamics.

3 Method & Problem Formulation

As this work tackles policy neural network weights generation, we take inspiration from the method employed in [22], where it was demonstrated that latent diffusion can be utilized to learn the distribution of policy parameters for humanoid locomotion. A key limitation of this method was the reliance on often unavailable policy datasets. This work does not require a dataset of policy parameters but trains on a (more commonly available) trajectory dataset.

LWD employs a two-step training process. First, a variational autoencoder (VAE) with weak KL-regularization encodes trajectories into a latent space that can be decoded into policy weights. The decoder is a conditioned hypernetwork that takes a policy neural network without its parameters (weights) as input and populates it with the desired weights based on the conditioning. The generated policy parameters are also optimized for trajectory tracking with a world model. Next, a diffusion model learns the distribution of this latent space, enabling policy sampling from the learned distribution (see Figure 2).

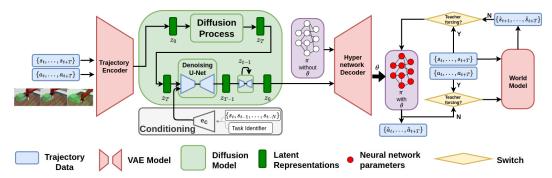


Figure 2: **LWD**: We first pre-train a VAE and World Model, which variationally encodes trajectories to a latent space and then decodes it as policy parameters, the policy parameters are optimized for both behavior cloning and trajectory tracking in a co-trained world model. Next, we train a conditional latent diffusion model to learn this latent distribution. When teacher forcing is enabled, the world model parameters are optimized, when disabled, the world model is used to optimize the VAE.

Compared to [22], which encodes policy parameters and employs a graph hypernetwork with a MSE loss on parameter reconstruction, our approach differs as it: (1) encodes trajectories as opposed to parameters, into latent space (2) uses a simple hypernetwork, (3) applies a behavior cloning loss (detailed in Section 3.1 & Section 3.2) on the generated policy, and (4) learns a world model for predicting observations given the action in an environment. Below we discuss the problem formulation and derivation.

3.1 Latent Policy Representation

We begin by formulating our approach for unconditional policy generation. Assume a distribution over stochastic policies, where variability reflects behavioral diversity. Each policy is parameterized by θ , with $\pi(\cdot,\theta)$ denoting a sampled policy and $p(\theta)$ the parameter distribution. Sampling a policy corresponds to drawing $\theta \sim p(\theta)$. When a policy interacts with the environment, it gives us a trajectory $\tau = \{s_t, a_t\}_{t=0}^T$. We assume multiple such trajectories are collected by repeatedly sampling θ and executing the corresponding policy. This enables a heterogeneous dataset, e.g., from humans or expert agents. For a given θ , actions are noisy: $a_t \sim \mathcal{N}(\pi(s_t,\theta),\sigma^2)$.

Our objective is to recover the distribution $p(\theta)$ that generated the trajectory dataset. We posit a latent variable z capturing behavioral modes, and assume conditional independence: $p(\tau \mid z, \theta) = p(\tau \mid \theta)$. Given trajectory data, we maximize the likelihood $\log p(\tau)$. To do so, we derive a modified Evidence Lower Bound (mELBO) that incorporates $p(\theta)$ (see below). This differs from the standard ELBO used in VAEs.

$$\log p(\tau) = \log \int \int p(\tau, \theta, z) \, dz \, d\theta \quad \text{(Introduce policy parameter } \theta \text{ and latent variable } z)$$

$$= \log \int \int p(\tau \mid z, \theta) p(\theta \mid z) p(z) \, dz \, d\theta \quad \text{(Apply the chain rule)}$$

$$= \log \int \int \frac{p(\tau \mid z, \theta) p(\theta \mid z) p(z)}{q(z \mid \tau)} q(z \mid \tau) \, dz \, d\theta \qquad \text{(1a)}$$

(Introduce a variational distribution $q(z \mid \tau)$, approximating the true posterior $p(z \mid \tau)$)

$$= \log \int \mathbb{E}_{p(\theta|z)} \left[\frac{p(\tau \mid z, \theta)p(z)}{q(z \mid \tau)} q(z \mid \tau) \right] dz \tag{1b}$$

$$\geq \mathbb{E}_{q(z|\tau)}\left[\log\left(\frac{\mathbb{E}_{p(\theta|z)}\left[p(\tau\mid z,\theta)\right]p(z)}{q(z\mid \tau)}\right)\right] \quad \text{(Jensen's inequality)}$$

$$= \mathbb{E}_{q(z\mid\tau)} \left[\log \left(\mathbb{E}_{p(\theta\mid z)} \left[p(\tau\mid z,\theta) \right] \right) \right] - \mathbb{E}_{q(z\mid\tau)} \left[\log \left(q(z\mid\tau) \right) - \log \left(p(z) \right) \right] \tag{1c}$$

$$=\mathbb{E}_{q(z\mid\tau)}\left[\log\left(\mathbb{E}_{p(\theta\mid z)}\left[p(\tau\mid\theta)\right]\right)\right]-\mathrm{KL}(q(z\mid\tau)\parallel p(z))\quad\text{(cond. independence)}\quad\text{(1d)}$$

$$\geq \mathbb{E}_{q(z\mid \tau)}\left[\mathbb{E}_{p(\theta\mid z)}\left[\log\left(p(\tau\mid \theta)\right)\right]\right] - \mathrm{KL}(q(z\mid \tau)\parallel p(z)) \quad \text{(Jensen's inequality)} \quad \text{(1e)}$$

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., $p(\tau \mid \theta)$) is given by:

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

 $\log p(s_1, a_1, \dots, s_T, a_T \mid \theta) = \log p(s_1) + \log p(a_1 \mid s_1, \theta)$

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

Substituting 2b in 1e:

$$\log p(\tau) \ge \mathbb{E}_{q(z|\tau)} \left[\mathbb{E}_{p(\theta|z)} \left[\log \left(p(\tau \mid \theta) \right) \right] \right] - \text{KL}(q(z \mid \tau) \parallel p(z))$$

$$= \mathbb{E}_{q(z|\tau)} \left[\mathbb{E}_{p(\theta|z)} \left[\sum_{t=1}^{T} \log p(a_t \mid s_t, \theta) + \sum_{t=2}^{T} \log p(s_t \mid s_{t-1}, a_{t-1}) \right] \right]$$

$$- \text{KL}(q(z \mid \tau) \parallel p(z)) + A$$
(3)

A consists of $\log p(s_1)$, and since this cannot be subject to maximization, we shall ignore it. Therefore, our modified ELBO is:

$$\mathbb{E}_{q(z|\tau)} \left[\mathbb{E}_{p(\theta|z)} \left[\sum_{t=1}^{T} \underbrace{\log p(a_t \mid s_t, \theta)}_{Behavior\ Cloning} + \sum_{t=2}^{T} \underbrace{\log p(s_t \mid s_{t-1}, a_{t-1})}_{World\ Model} \right] \right] - \underbrace{\text{KL}(q(z \mid \tau) \parallel p(z))}_{KL\ Regularizer}$$
(4)

3.2 Loss function for World Model Augmented Variational Autoencoder for Policies

Since we now have a modified ELBO objective, we shall now try to approximate its components with a variational autoencoder and a world model. Let ϕ_{enc} be the parameters of the VAE encoder that variationally maps trajectories to z, ϕ_{dec} be the parameters of the VAE decoder, and ϕ_{wm} be the world model parameters. We assume the latent z is distributed with mean zero and unit variance. We construct the VAE decoder to approximate $p(\theta \mid z)$ with $p_{\phi_{dec}}(\theta \mid z)$. Considering $a_t \sim \mathcal{N}(\pi(s_t, \theta), \sigma^2)$, and $\tau_k = \{s_t^k, a_t^k\}_{t=1}^T$, we derive our VAE loss function as:

$$\mathcal{L}_{BC} = -\sum_{t=1}^{T} \mathbb{E}_{q_{\phi_{enc}}(z|\tau_{k})} \left[(a_{t}^{k} - \pi(s_{t}^{k}, f_{\phi_{dec}}(z)))^{2} \right]$$

$$\mathcal{L}_{RO} = -\sum_{t=2}^{T} \mathbb{E}_{q_{\phi_{enc}}(z|\tau_{k})} \left[\text{KL} \left(p_{\phi_{wm}}(s_{t} \mid s_{t-1}^{k}, \pi(s_{t-1}^{k}, f_{\phi_{dec}}(z))) \parallel p_{\phi_{wm}}(s_{t} \mid s_{t-1}^{k}, a_{t-1}^{k}) \right) \right]$$

$$\mathcal{L}_{TF} = -\sum_{t=2}^{T} (s_{t}^{k} - \hat{s}_{t}^{k})^{2} \qquad \mathcal{L}_{KL} = -\beta_{kl} \sum_{i=1}^{\dim(z)} \left(\sigma_{e_{i}}^{2} + \mu_{e_{i}}^{2} - 1 - \log \sigma_{e_{i}}^{2} \right)$$

$$\mathcal{L} \left(\left\{ s_{t}^{k}, a_{t}^{k} \right\}_{t=1}^{T} \mid \phi_{enc}, \phi_{dec}, \phi_{wm} \right) = \mathcal{L}_{BC} + \mathcal{L}_{RO} + \mathcal{L}_{TF} + \mathcal{L}_{KL}$$
(5)

where, \mathcal{L}_{BC} is the behavior cloning loss to train the policy decoder, \mathcal{L}_{RO} is the rollout loss to correct the decoded policy's actions using the world model, \mathcal{L}_{TF} is the teacher forcing loss to train the world model, and \mathcal{L}_{KL} is the KL loss to regularize the latent space. θ is obtained from the hypernetwork decoder $f_{\phi_{dec}}(z)$. $(\mu_e, \sigma_e) = f_{\phi_{enc}}(\{s_t^k, a_t^k\}_{t=1}^T), z \sim \mathcal{N}(\mu_e, \sigma_e), \hat{s}_t^k \sim p_{\phi_{wm}}(s_t^k \mid s_{t-1}^k, a_{t-1}^k)$ and β_{kl} is the regularization weight. The complete derivation is shown in Section A. Since the decoder in the VAE outputs the parameter of a secondary network, we shall use a conditional hypernetwork, specifically the model developed for continual learning by [53]. For computational stability, we shall use \mathcal{L}_{BC} , \mathcal{L}_{RO} and \mathcal{L}_{KL} to optimize the VAE (encoder and decoder parameters) and \mathcal{L}_{TF} to train the world model parameters. With the teacher forcing objective we get a reliable world model that we can then use in the rollout objective. This is similar to procedures followed in [1, 44, 27].

In practice, we see that approximating $p(z) = \mathcal{N}(0, I)$ is suboptimal, and therefore we set β_{kl} to a very small number $\sim (10^{-10}, 10^{-6})$. After training the VAE to maximize the objective provided

in Equation (5) with this β_{kl} , we have access to this latent space z and can train a diffusion model to learn its distribution p(z). We can condition the latent denoising process on the current state and/or the task identifier c of the policy required. Therefore the model shall be approximating $p_{\phi_{dif}}(z_{t-1} \mid z_t, c)$. After denoising for a given state and task identifier, we can convert the denoised latent to the required policy. Therefore, to sample from $p(\theta)$, first sample z using the trained diffusion model $z \sim p_{\phi_{dif}}(z_0)$, and then apply the deterministic function $f_{\phi_{dec}}$ to the sampled z. Note that to sample policies during inference, we do not need to encode trajectories; rather, we need to sample a latent using the diffusion model and use the hypernetwork decoder of a pre-trained VAE to decode a policy from it.

4 Experiments

We run three sets of experiments. In the first set, we evaluate the validity of our main contributions. We compare LWD with action trajectory generation methods with respect to 1) Longer Action Horizons and Environment Perturbations, where experiments are performed on the PushT task [9] and the Lift and Can Robomimic tasks [36], while varying these parameters, and 2) Lower inference costs, where experiments are performed on 10 tasks from the Metaworld [58] suite of tasks, to show LWD requires fewer parameters during inference while maintaining multi-task performance. The task descriptions are provided in Section B. We choose a multi-task experiment here as the model capacity required for solving multiple tasks generally increases with the number of tasks. In the second set, we perform ablations over three components of our method: 1) Diffusion model architecture, Section E; 2) VAE decoder size, Section F; 3) KL coefficient for the VAE, Section D. In the final set, we analyze the behavior distribution modeled by our latent space. These can be found in Section G, and Section H.

We focus on demonstrating results in state-based low-dimensional observation spaces. Our generated policies are Multi-Layer Perceptrons (MLP) with 2 hidden layers with 256 neurons each. In the VAE, the encoder is a sequential network that flattens the trajectory and compresses it to a lowdimensional latent space, and the decoder is a conditional hypernetwork[12]. The details of the VAE implementation are provided in Section J.2 and Section J.3. For the world model, since we use low-dimensional observation spaces, we use a simple MLP with 2 hidden layers with 1024 neurons each to map the history of observations and actions to the next observation. For stability, we use \mathcal{L}_{RO} only after 10 epochs of training. This warm-starts the world model before we use it to optimize the policy generator. For all experiments, the latent space is \mathbb{R}^{256} and the learning rate is 10^{-4} with the Adam optimizer. For the diffusion model, we use the DDPM Scheduler for denoising. Inspired by [9], we conducted an ablation between two diffusion architectures: a ConditionalUnet1D network and a Transformer architecture. Based on the results are shown in Section E, We chose the ConditionalUnet1D model for all experiments in the paper. Just as [9], we condition the diffusion model with FiLM layers, and also use the Exponential Moving Average [21] of parameter weights (commonly used in DDPM) for stability. All results presented in this work are obtained from running experiments over three seeds, and the compute resources are described in Section J.9

4.1 Empirical Evaluation of Contributions

4.1.1 Longer Action Horizons & Robustness to Perturbations

We first evaluate our method on the PushT task [9], a standard benchmark for diffusion-based trajectory generation in manipulation. The goal is to align a 'T' block with a target position and orientation on a 2D surface. Observations consist of the end-effector's position and the block's position and orientation. Actions specify the end-effector's target position at each time step. Task performance (measured by success rate) is defined as the maximum overlap between the actual and desired block poses during a rollout. We test under different action horizons and varying levels of environment perturbation, simulated via an adversarial agent that randomly displaces the 'T' block.

For the LWD model, we first train a VAE to encode trajectory snippets (of length equal to the action horizon) into latents representing locally optimal policies. These policies are optimized with a co-trained world model. A conditional latent diffusion model, given the current state, then generates a latent that the VAE decoder transforms into a locally optimal policy for the next action horizon. The inference process is illustrated in Figure 1. We train two variants of LWD, with (LWD w/ WM) and without (LWD w/o WM) the world model (i.e., we train LWD with just $\mathcal{L}_{BC} + \mathcal{L}_{KL}$).

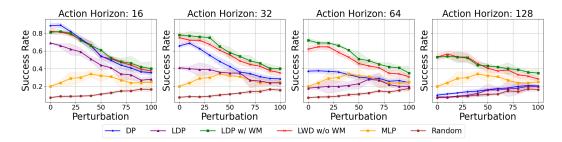


Figure 3: Longer action horizons and robustness to perturbations on PushT: Performance of LWD and baselines on the PushT task as we vary the action horizon and environment perturbations.

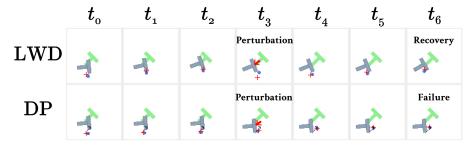


Figure 4: **Visualization of Perturbation**: When an adversarial perturbation is applied, we see that LWD's generated closed-loop policy successfully adapts to the change.

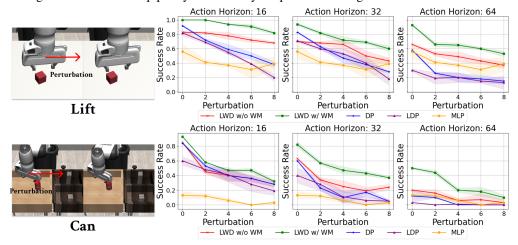


Figure 5: Longer action horizons and robustness to perturbations on Robomimic tasks: Performance of LWD and DP as we vary the action horizon and environment perturbations.

As baselines for this experiment, we compare the proposed LWD variants against four alternatives: 1) a **Diffusion Policy (DP)** model that generates open-loop action trajectories for a fixed action horizon; 2) a **Latent Diffusion Policy (LDP)** model, which is structurally similar to LWD but decodes the latent representation into an action trajectory rather than a closed-loop policy; 3) a **Multilayer Perceptron (MLP)** policy, which shares the same architecture as the policy network generated by LWD and serves to isolate the impact of diffusion modeling; 4) a **Random Policy**, which provides a lower-bound performance reference. For a fair comparison, all diffusion-based models (LWD, DP, and LDP) use the same model size and hyperparameters, corresponding to the medium configuration described in Section J.4 and Section J.7. LDP uses a VAE decoder, implemented as an MLP with two hidden layers of 256 neurons each, to output an action chunk of the same length as the action horizon.

All models are evaluated across 50 uniquely seeded environment instances, with each evaluation repeated 10 times, across 3 training seeds. Figure 3 illustrates the impact of perturbation magnitudes and action horizons on success rates across all baselines. Perturbations refer to random displacements applied to the T block, occurring at randomly selected time steps with 10% probability. A sample rollout with a perturbation magnitude of 50 is shown in Figure 4.

While DP demonstrates comparable performance to both LWD variants at an action horizon of 16 with minimal perturbations, LWD exhibits superior robustness as the action horizon increases. This enhanced robustness of LWD with the world model becomes more pronounced in the presence of larger perturbations. Specifically, at longer action horizons such as 128, LWD w/ WM maintains a significantly higher success rate compared to DP across all perturbation levels. The MLP generally underperforms compared to both LWD variants and DP, highlighting the benefits of diffusion-based approaches for this task. LDP has a lower success rate than LWD, indicating that generating a closed-loop policy is more important than learning the latent representation space. The relatively lower sensitivity to perturbations at an action horizon of 16 for both policies can be attributed to the more frequent action trajectory queries inherent in DP at shorter horizons (i.e. smaller action chunks), effectively approximating a more closed-loop control strategy.

Finally, we measured the total time required to successfully complete the PushT task and see that for all levels of perturbation, for a given success rate, LWD executes a rollout quicker (in wall-clock time) than the DP model. See Section I for more details.

We also ran experiments on the Robomimic [36] Lift and Can tasks, using the same hyperparameters as the PushT experiment, the same task settings, and the mh demonstration data from [9]. To simulate perturbations, we add random translation and rotation vectors to the end effector, applied 10% of the time. Figure 5 shows the performance of the LWD variants and baselines under these perturbations across different action horizons. The x-axis corresponds to perturbation magnitude. Similar to PushT, LWD outperforms DP for longer horizons and is more robust to perturbations.

4.1.2 Low Inference Cost

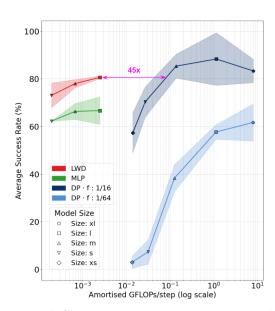


Figure 6: Success rate vs. average compute of LWD, DP, and MLP policies on 10 Metaworld tasks for various model sizes. The x-axis shows the number of GFLOPS/step for each policy on a log scale. performs ~ 45 times fewer computations than a DP policy with comparable performance.

We will now look at the next contribution, namely, lower inference cost compared to methods that diffuse action trajectories instead of policies. When training a single policy on multiple tasks, it is known that a larger model capacity is needed. This is detrimental in robotics applications as this increases control latency. We train a task-conditioned LWD model and show that the cost of task generalization is borne by the latent diffusion model, while the generated execution policy remains small. Because LWD generates a smaller policy, the runtime compute required for inference is much smaller compared to the SOTA methods. For simplicity, in this section, we use the LWD w/o WM variant.

We experiment on 10 tasks of the Metaworld benchmark, the details of which are in Section B. We set the action horizon to the length of the entire trajectory for LWD to generate policies that shall work for the entire duration of the rollout, where at each time step, the generated MLPs shall predict instantaneous control. We experimented over three sizes of the generated MLP policy: 128, 256, and 512 neurons per layer, each having 2 hidden layers. We also train 10 DP models, spread over a grid of 5 different sizes (xs, s, m, 1, xl) and 2 action horizons: 32

and 128. Each DP model is run at an inference frequency of half the action horizon. We provide the details of the DP model in Section J.1. Finally, we also train 3 MLP models with 128, 256, and 512 neurons per layer, to compare the performance of LWD with a standard MLP policy.

Note that LWD uses a fixed action horizon equal to the full episode length (500 steps), whereas the DP model uses a variable horizon. The LWD inference process is illustrated on the right-hand side of Figure 1. All baseline models receive the task identifier as part of the state input. Each model is trained with 3 random seeds, and evaluated across 10 tasks, with 16 rollouts per task. Figure 6 presents the results of this evaluation. In the plot, the x-axis represents average per-step inference

compute (in GFLOPs), and the y-axis indicates the overall success rate across tasks. For DP models, achieving high success rates requires increasing model size or denoising frequency (i.e., predicting shorter action chunks), both of which raise computational cost. In contrast, LWD generates a simpler, more efficient controller, requiring significantly less compute. The best-performing LWD model achieves an 81% success rate with $\sim45\times$ fewer inference operations than the closest-performing DP model. Interestingly, the MLP baseline also performs well, and is comparable in efficiency to LWD, but still lags in performance. We attribute this to the unimodal nature of this dataset, as MLPs struggled with the multimodal PushT task in the previous section.

4.2 Behavior Analysis

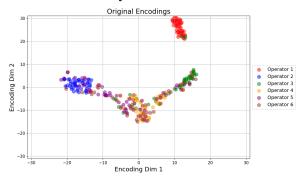


Figure 7: Behavior distribution: Robomimic Lift task

LWD models trajectory data from a distribution of policies, exposing this distribution through its latent space. On the Robomimic Lift task with the MH dataset (300 trajectories from 6 operators of varied proficiency: 2 "worse," 2 "okay," and 2 "better."), LWD encoded entire demonstration trajectories. A 2D t-SNE plot revealed clusters aligned with operator identity, despite LWD receiving no explicit operator labels. This shows LWD can cluster behaviors and potentially filter unwanted ones, a capability further studied in Section G and Section H.

5 Limitations and Future Work

Latent Weight Diffusion (LWD) is a promising framework for policy generation, but Diffusion Policy (DP) performs better in short-horizon, low-perturbation settings. This gap likely stems from VAE approximation errors and LWD's added complexity, which increases training cost, particularly the RAM demands of loading full observation sequences.

DP's strength is its compatibility with foundation models as action heads [40], while integrating LWD into such architectures remains an open question. For image-based observations, LWD requires CNN-based policies. Although hypernetworks can generate CNN and ViT weights, adapting these within LWD is challenging. Our early results show promise in using pre-trained encoders and applying LWD to image embeddings. Integrating LWD with vision-based world models is also open, but appears feasible with approaches in [44] and [27].

Future work could improve LWD's VAE decoder through chunked deconvolutional hypernetworks [53], enabling more efficient decoding. Extending LWD to Transformer or ViT policies is another direction, especially for sequential or visual tasks [11]. Finally, warm-starting with prior latents [9] may further boost performance by providing richer priors.

6 Conclusion

We introduce Latent Weight Diffusion (LWD), a novel framework for learning a distribution over policies from diverse demonstration trajectories. LWD models behavioral diversity via latent diffusion, a world model, and uses a hypernetwork decoder to generate policy weights, enabling closed-loop control directly from sampled latents. Our evaluation highlights two key strengths of LWD: robustness and computational efficiency. Compared to Diffusion Policy, LWD delivers more reliable performance in environments with long action horizons and perturbations, while reducing inference costs, especially in multi-task settings.

References

- [1] Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025.
- [2] Sumeet Batra, Bryon Tjanaka, Matthew C Fontaine, Aleksei Petrenko, Stefanos Nikolaidis, and Gaurav Sukhatme. Proximal policy gradient arborescence for quality diversity reinforcement learning. *arXiv preprint arXiv:2305.13795*, 2023.
- [3] Luca Bertinetto, João F. Henriques, Jack Valmadre, Philip Torr, and Andrea Vedaldi. Learning feed-forward one-shot learners. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 29. Curran Associates, Inc., 2016.
- [4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alex Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *arXiv preprint arXiv:2307.15818*, 2023.
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In arXiv preprint arXiv:2212.06817, 2022.
- [6] Joao Carvalho, An T Le, Mark Baierl, Dorothea Koert, and Jan Peters. Motion planning diffusion: Learning and planning of robot motions with diffusion models. in 2023 ieee. In RSJ International Conference on Intelligent Robots and Systems (IROS), pages 1916–1923.
- [7] Shang-Fu Chen, Hsiang-Chun Wang, Ming-Hao Hsu, Chun-Mao Lai, and Shao-Hua Sun. Diffusion model-augmented behavioral cloning. In *International Conference on Machine Learning*, pages 7486–7510. PMLR, 2024.
- [8] Yixuan Chen, Hao Zhang, and Jian Liu. Time-aware world model for adaptive prediction and control. *arXiv preprint arXiv:2506.08441*, 2025.
- [9] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 2024.
- [10] Open X-Embodiment Collaboration, Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang,

Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan 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 Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. https://arxiv.org/abs/2310.08864, 2023.

- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *CoRR*, abs/2010.11929, 2020.
- [12] Benjamin Ehret, Christian Henning, Maria R. Cervera, Alexander Meulemans, Johannes von Oswald, and Benjamin F. Grewe. Continual learning in recurrent neural networks. In *International Conference on Learning Representations*, 2021.
- [13] Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pages 158–168. PMLR, 08–11 Nov 2022.
- [14] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning, 2020.

- [15] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. arXiv preprint arXiv:1609.09106, 2016.
- [16] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems 31*, pages 2451–2463. Curran Associates, Inc., 2018. https://worldmodels.github.io.
- [17] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- [18] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. *arXiv preprint arXiv:1811.04551*, 2019.
- [19] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Mastering atari with discrete world models. arXiv preprint arXiv:2010.02193, 2020.
- [20] Danijar Hafner, Jurgis Pašukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- [21] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- [22] Shashank Hegde, Sumeet Batra, KR Zentner, and Gaurav Sukhatme. Generating behaviorally diverse policies with latent diffusion models. *Advances in Neural Information Processing Systems*, 36:7541–7554, 2023.
- [23] Shashank Hegde, Zhehui Huang, and Gaurav S Sukhatme. Hyperppo: A scalable method for finding small policies for robotic control. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 10821–10828. IEEE, 2024.
- [24] Negin Heravi, Ayzaan Wahid, Corey Lynch, Peter R. Florence, Travis Armstrong, Jonathan Tompson, Pierre Sermanet, Jeannette Bohg, and Debidatta Dwibedi. Visuomotor control in multi-object scenes using object-aware representations. 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 9515–9522, 2022.
- [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc' Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- [26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [27] Anthony Hu, Zak Murez, Nikhil Mohan, Sofia Dudas, Jeffrey Hawke, Vijay Badrinarayanan, Roberto Cipolla, and Alex Kendall. Model-based imitation learning for urban driving. *arXiv* preprint arXiv:2210.07729, 2022.
- [28] Xiaoyu Huang, Yufeng Chi, Ruofeng Wang, Zhongyu Li, Xue Bin Peng, Sophia Shao, Borivoje Nikolic, and Koushil Sreenath. Diffuseloco: Real-time legged locomotion control with diffusion from offline datasets. *arXiv preprint arXiv:2404.19264*, 2024.
- [29] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [30] Deqian Kong, Dehong Xu, Minglu Zhao, Bo Pang, Jianwen Xie, Andrew Lizarraga, Yuhao Huang, Sirui Xie, and Ying Nian Wu. Latent plan transformer for trajectory abstraction: Planning as latent space inference. *Advances in Neural Information Processing Systems*, 37:123379–123401, 2024.

- [31] Alex X Lee, Anusha Nagabandi, Pieter Abbeel, and Sergey Levine. Stochastic latent actor-critic: Deep reinforcement learning with a latent variable model. *arXiv preprint arXiv:1907.00953*, 2019.
- [32] Yongyuan Liang, Tingqiang Xu, Kaizhe Hu, Guangqi Jiang, Furong Huang, and Huazhe Xu. Make-an-agent: A generalizable policy network generator with behavior-prompted diffusion. *arXiv preprint arXiv:2407.10973*, 2024.
- [33] Zhixuan Liang, Yao Mu, Hengbo Ma, Masayoshi Tomizuka, Mingyu Ding, and Ping Luo. Skilldiffuser: Interpretable hierarchical planning via skill abstractions in diffusion-based task execution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16467–16476, 2024.
- [34] Justin Lovelace, Varsha Kishore, Chao Wan, Eliot Shekhtman, and Kilian Q Weinberger. Latent diffusion for language generation. Advances in Neural Information Processing Systems, 36, 2024.
- [35] Yunhao Luo, Chen Sun, Joshua B Tenenbaum, and Yilun Du. Potential based diffusion motion planning. *arXiv preprint arXiv:2407.06169*, 2024.
- [36] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *arXiv preprint arXiv:2108.03298*, 2021.
- [37] Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models. In *International Conference on Learning Representations*, 2023.
- [38] Vincent Micheli, Eloi Alonso, and François Fleuret. Storm: Efficient stochastic transformer based world models for reinforcement learning. *arXiv preprint arXiv:2310.09615*, 2024.
- [39] Utkarsh Aashu Mishra, Shangjie Xue, Yongxin Chen, and Danfei Xu. Generative skill chaining: Long-horizon skill planning with diffusion models. In *Conference on Robot Learning*, 2023.
- [40] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Lawrence Yunliang Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- [41] William Peebles, Ilija Radosavovic, Tim Brooks, Alexei A Efros, and Jitendra Malik. Learning to learn with generative models of neural network checkpoints. *arXiv preprint arXiv:2209.12892*, 2022.
- [42] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings* of the IEEE/CVF international conference on computer vision, pages 4195–4205, 2023.
- [43] Xue Peng, Erwin Coumans, Tingnan Zhang, Tsang-Wei Lee, Jie Tan, and Sergey Levine. Learning agile robotic locomotion skills by imitating animals. 07 2020.
- [44] Alexander Popov, Alperen Degirmenci, David Wehr, Shashank Hegde, Ryan Oldja, Alexey Kamenev, Bertrand Douillard, David Nistér, Urs Muller, Ruchi Bhargava, et al. Mitigating covariate shift in imitation learning for autonomous vehicles using latent space generative world models. *arXiv preprint arXiv:2409.16663*, 2024.
- [45] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2018.
- [46] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned imitation learning using score-based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.

- [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pages 10674–10685. IEEE, 2022.
- [49] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.
- [50] Aviv Tamar, Yi Wu, Garrett Thomas, Sergey Levine, and Pieter Abbeel. Value iteration networks. In *Advances in Neural Information Processing Systems*, pages 2154–2162, 2016.
- [51] Wenhui Tan, Bei Liu, Junbo Zhang, Ruihua Song, and Jianlong Fu. Multi-task manipulation policy modeling with visuomotor latent diffusion. *ArXiv*, abs/2403.07312, 2024.
- [52] Julen Urain, Niklas Funk, Jan Peters, and Georgia Chalvatzaki. Se(3)-diffusionfields: Learning smooth cost functions for joint grasp and motion optimization through diffusion. 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 5923–5930, 2022.
- [53] Johannes von Oswald, Christian Henning, Benjamin F. Grewe, and João Sacramento. Continual learning with hypernetworks. In *International Conference on Learning Representations*, 2020.
- [54] Kai Wang, Dongwen Tang, Boya Zeng, Yida Yin, Zhaopan Xu, Yukun Zhou, Zelin Zang, Trevor Darrell, Zhuang Liu, and Yang You. Neural network diffusion. arXiv preprint arXiv:2402.13144, 2024.
- [55] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in Neural Information Processing Systems*, pages 2746–2754, 2015.
- [56] Philipp Wu, Alejandro Escontrela, Danijar Hafner, Ken Goldberg, and Pieter Abbeel. Day-dreamer: World models for physical robot learning. arXiv preprint arXiv:2206.14176, 2022.
- [57] Zhou Xian and Nikolaos Gkanatsios. Chaineddiffuser: Unifying trajectory diffusion and keypose prediction for robotic manipulation. In *Conference on Robot Learning/Proceedings of Machine Learning Research*. Proceedings of Machine Learning Research, 2023.
- [58] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pages 1094–1100. PMLR, 2020.
- [59] Hao Zhang, Zhihan Xu, Jian Liu, and Qingzhao Wang. Generalized and efficient planning with scalable latent world models. *arXiv preprint arXiv:2406.10667*, 2024.
- [60] Weidong Zhang, Jian Liu, Lihe Xia, Qingzhao Wang, and Hongming Zhou. Safedreamer: Safe reinforcement learning with world models. *arXiv preprint arXiv:2307.07176*, 2023.
- [61] Ruijie Zheng, Ching-An Cheng, Hal Daumé Iii, Furong Huang, and Andrey Kolobov. Prise: Llm-style sequence compression for learning temporal action abstractions in control. In *International Conference on Machine Learning*, pages 61267–61286. PMLR, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Yes, the paper presents a novel approach to generating reactive policies using latent weight diffusion, which is articulated in both the intro and the abstract sections.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitations and Future Work are discussed in Section 5.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We have a proof in Section 3.1 and Section A that is complete and correct to the best of the authors' knowledge.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The authors have provided full information about the method, the setup for experiments, and the hyperparameters (in the appendices) used to reproduce these results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We will include the code and data in the camera-ready version.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Yes, these details are provided in the paper, with more details (such as hyperparameters) provided in the appendices.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All experiemnts were trained with three seeds. For PushT tasks, we ran 50 evaluations for each seed. For MetaWorld, we ran 16 evaluations per seed per task. For Robomimic, we ran 32 evaluations per seed per task. Each plot has error bars.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: This is mentioned in Section J.9. We have even more precise calculations for evaluation compute resources measured in Giga FLOPS when comparing compute with the baseline (DP).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The authors have adhered to the NeurIPS Code of Ethics.

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: A societal impact section is mentioned in Section K.

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper uses datasets for solving robotic tasks, which are not high risk for misuse.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The authors have credited the papers where the data has been used to train the models. The licenses were permissive licenses (Apache 2.0 license for D4RL, MIT license for Diffusion Policy) to use for research purposes. The authors have also cited the original papers that produced the code package or dataset.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The assets here are the code and the models. The authors will release these for the camera-ready version.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No human subjects were involved in this work.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This work does not involve human subjects, and did not require IRB approval.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: This paper does not involve LLMs as any component of the core methods in this research.

A Appendix – VAE loss derivation

Since $a_t \sim \mathcal{N}(\pi(s_t, \theta), \sigma^2)$:

$$p(a_t \mid s_t, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(a_t - \pi(s_t, \theta))^2}{2\sigma^2}\right)$$
 (6)

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

$$\mathcal{L}\left(\tau_{k} \mid \phi_{enc}, \phi_{dec}, \phi_{wm}\right) \\
= \sum_{t=1}^{T} \mathbb{E}_{q_{\phi_{enc}}(z \mid \tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta \mid z)} \left[\log p\left(a_{t}^{k} \mid s_{t}^{k}, \theta\right) \right] \right] \\
+ \sum_{t=2}^{T} \mathbb{E}_{q_{\phi_{enc}}(z \mid \tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta \mid z)} \left[\log p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, a_{t-1}^{k}) \right] \right] \\
- \text{KL}\left(q_{\phi_{enc}}\left(z \mid \tau_{k}\right) \parallel p(z)\right) \tag{7}$$

Consider the second term in the above equation. On maximization $a_{t-1}^k = \pi(s_{t-1}^k, \theta)$, and because the inner quantity is a constant w.r.t. s_t we can add a harmless expectation $\mathbb{E}_{s_t \sim \pi}[.]$ (i.e., states visited by the estimated policy, not necessarily those in the dataset), therefore it becomes:

$$\mathbb{E}_{q_{\phi_{enc}}(z|\tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta|z)} \left[\mathbb{E}_{s_{t} \sim \pi} \left[\log p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, \pi(s_{t-1}^{k}, \theta)) \right] \right] \right] \\
= \mathbb{E}_{q_{\phi_{enc}}(z|\tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta|z)} \left[\mathbb{E}_{s_{t} \sim \pi} \left[\log \frac{p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, \pi(s_{t-1}^{k}, \theta))}{p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, a_{t-1}^{k})} \right] \right] \right] \\
+ \mathbb{E}_{q_{\phi_{enc}}(z|\tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta|z)} \left[\mathbb{E}_{s_{t} \sim \pi} \left[\log p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, a_{t-1}^{k}) \right] \right] \right] \tag{8}$$

We can now substitute in the KL term, and drop the expectation in the last term (since the inner terms only depend on s_{t-1}^k and not $s_t \sim \pi$, θ , or z. Therefore, the loss now becomes:

$$\mathcal{L}\left(\tau_{k} \mid \phi_{enc}, \phi_{dec}, \phi_{wm}\right) \\
= C - \frac{1}{2\sigma^{2}} \sum_{t=1}^{T} \mathbb{E}_{q_{\phi_{enc}}(z \mid \tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta \mid z)} \left[(a_{t}^{k} - \pi(s_{t}^{k}, \theta))^{2} \right] \right] \\
- \sum_{t=2}^{T} \mathbb{E}_{q_{\phi_{enc}}(z \mid \tau_{k})} \left[\mathbb{E}_{p_{\phi_{dec}}(\theta \mid z)} \left[\text{KL} \left(p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, \pi(s_{t-1}^{k}, \theta)) \parallel p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, a_{t-1}^{k}) \right) \right] \right] \\
- \sum_{t=2}^{T} \log p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, a_{t-1}^{k}) \\
- \text{KL} \left(q_{\phi_{enc}}(z \mid \tau_{k}) \parallel p(z) \right) \tag{9}$$

For computational stability, we construct our decoder to be a deterministic function $f_{\phi_{dec}}$, i.e., $p_{\phi_{dec}}(\theta \mid z)$ becomes $\delta(\theta - f_{\phi_{dec}}(z))$. Further, if we have a trained world model, we can approximate s_t^k with s_t (i.e., direct model output samples) in the second term. This is done so that we can optimize the world model and policy correction separately with the teacher forcing and rollout objectives

(similar to that followed in [1]. Therefore:

$$\begin{split} &\mathcal{L}\left(\tau_{k} \mid \phi_{enc}, \phi_{dec}, \phi_{wm}\right) \\ &= C - \frac{1}{2\sigma^{2}} \sum_{t=1}^{T} \mathbb{E}_{q_{\phi_{enc}}(z \mid \tau_{k})} \left[(a_{t}^{k} - \pi(s_{t}^{k}, f_{\phi_{dec}}(z)))^{2} \right] \\ &- \sum_{t=2}^{T} \mathbb{E}_{q_{\phi_{enc}}(z \mid \tau_{k})} \left[\text{KL} \left(p_{\phi_{wm}}(s_{t} \mid s_{t-1}^{k}, \pi(s_{t-1}^{k}, f_{\phi_{dec}}(z))) \mid \mid p_{\phi_{wm}}(s_{t} \mid s_{t-1}^{k}, a_{t-1}^{k}) \right) \right] \\ &- \sum_{t=2}^{T} \log p_{\phi_{wm}}(s_{t}^{k} \mid s_{t-1}^{k}, a_{t-1}^{k}) \\ &- \text{KL} \left(q_{\phi_{enc}}(z \mid \tau_{k}) \mid \mid p(z) \right) \end{split}$$

Where C is a constant from the substitution. Enforcing $p(z) = \mathcal{N}(0, I)$, and ignoring constants, we get:

$$\mathcal{L}\left(\tau_{k} \mid \phi_{enc}, \phi_{dec}, \phi_{wm}\right) = \mathcal{L}_{BC} + \mathcal{L}_{RO} + \mathcal{L}_{TF} + \mathcal{L}_{KL} \tag{10}$$

$$\mathcal{L}_{BC} = -\sum_{t=1}^{T} \mathbb{E}_{q_{\phi_{enc}}(z|\tau_k)} \left[(a_t^k - \pi(s_t^k, f_{\phi_{dec}}(z)))^2 \right]$$
 (11)

$$\mathcal{L}_{RO} = -\sum_{t=2}^{T} \mathbb{E}_{q_{\phi_{enc}}(z|\tau_{k})} \left[\text{KL} \left(p_{\phi_{wm}}(s_{t} \mid s_{t-1}^{k}, \pi(s_{t-1}^{k}, f_{\phi_{dec}}(z))) \mid \mid p_{\phi_{wm}}(s_{t} \mid s_{t-1}^{k}, a_{t-1}^{k}) \right) \right]$$
(12)

$$\mathcal{L}_{TF} = -\sum_{t=2}^{T} (s_t^k - \hat{s}_t^k)^2$$
 (13)

$$\mathcal{L}_{KL} = -\beta_{kl} \sum_{i=1}^{\dim(z)} \left(\sigma_{e_i}^2 + \mu_{e_i}^2 - 1 - \log \sigma_{e_i}^2 \right)$$
 (14)

where, \mathcal{L}_{BC} is the behavior cloning loss to train the policy decoder, \mathcal{L}_{RO} is the rollout loss to correct the decoded policy's actions using the world model, \mathcal{L}_{TF} is the teacher forcing loss to train the world model, \mathcal{L}_{KL} is the KL loss to regularize the latent space, $(\mu_e, \sigma_e) = f_{\phi_{enc}}(\tau_k)$, $z \sim \mathcal{N}(\mu_e, \sigma_e)$, $\hat{s}_t^k \sim p_{\phi_{wm}}(s_t^k \mid s_{t-1}^k, a_{t-1}^k)$ and β_{kl} is the regularization weight.

B Appendix – Metaworld task descriptions

Task	Description
Window Open	Push and open a window. Randomize window positions
Door Open	Open a door with a revolving joint. Randomize door positions
Drawer Open	Open a drawer. Randomize drawer positions
Dial Turn	Rotate a dial 180 degrees. Randomize dial positions
Faucet Close	Rotate the faucet clockwise. Randomize faucet positions
Button Press	Press a button. Randomize button positions
Door Unlock	Unlock the door by rotating the lock clockwise. Randomize door positions
Handle Press	Press a handle down. Randomize the handle positions
Plate Slide	Slide a plate into a cabinet. Randomize the plate and cabinet positions
Reach	Reach a goal position. Randomize the goal positions

Table 1: Metaworld task descriptions and randomization settings

C Appendix – Effect of Trajectory snipping on Latent Representations

For most robotics use cases, it is impossible to train on long trajectories due to the computational limitations of working with large batches of long trajectories. In some cases, it may also be beneficial to generate locally optimum policies for shorter action horizons (as done for experiments presented in Section 4.1.1). Therefore, we analyze the effect of sampling smaller sections of trajectories from the dataset. After training a VAE for the D4RL half-cheetah dataset on three policies (expert, medium, and random), we encode all the trajectories in the mixed dataset to the latent space. We then perform Principal Component Analysis (PCA) on this set of latents and select the first two principal components. Figure 8a shows us a visualization of this latent space. We see that the VAE has learned to encode the three sets of trajectories to be well separable. Next, we run the same experiment, but now we sample trajectory snippets of length 100 from the dataset instead of the full-length (1000) trajectories, Figure 8b shows us the PCA on the encoded latents of these trajectory snippets. We see that the separability is now harder in the latent space. Surprisingly, we noticed that after training our VAE on the snippets, the decoded policies from randomly snipped trajectories were still faithfully behaving like their original policies. We believe that this is because the halfcheetah env is a cyclic locomotion task, and all trajectory snippets have enough information to indicate its source policy. More dimensions of the PCA are shown in Figure 9.

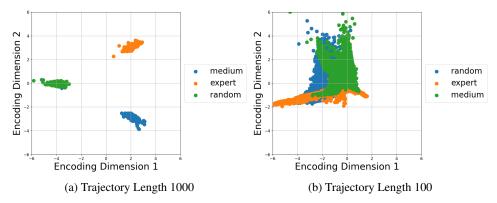


Figure 8: Effect of trajectory snipping in HalfCheetah. Top two principal components of the latent.

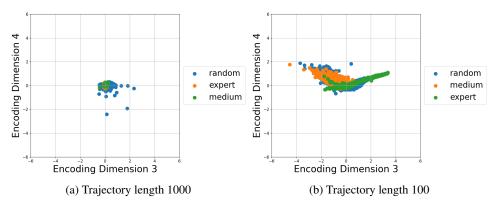


Figure 9: **Effect of trajectory snipping** in HalfCheetah. Top third and fourth principal components of the latent.

To validate this hypothesis, we analyze our method on trajectory snippets for non-cyclic tasks. We choose the MT10 suite of tasks in Metaworld [58] (note that these are different from the original 10 tasks discussed in the rest of the paper. We utilize the hand-crafted expert policy for each of the tasks in MT10 to collect trajectory data. For each task, we collect 1000 trajectories of length 500.

Figure 10a shows the principal components of the latents of the full trajectories in the dataset, and Figure 10b shows the same for the split trajectories. We can see that the separability of different tasks is much harder in this case. More dimensions of the PCA are shown in Figure 11b. Further, we

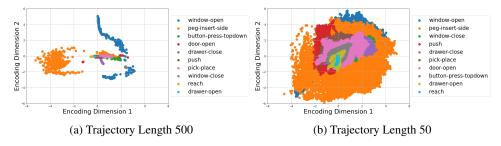


Figure 10: Effect of trajectory snipping in MT10. Top two principal components of the latent.

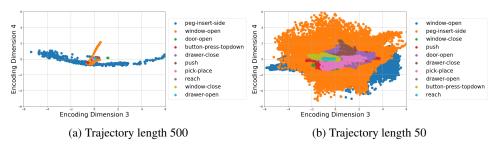


Figure 11: **Effect of trajectory snipping** in MT10. Top third and fourth principal components of the latent.

noticed that the decoded policies from the trajectory snippets did not perform as well as the original policies - for the same decoder size as the half cheetah task. This validates our hypothesis that the snippets are unable to reproduce the original policy for non-cyclic tasks. To have the same degree of behavior reconstruction as the half-cheetah tasks, we need a larger decoder model. This is discussed in Section F.

D Appendix - Ablation of the KL coefficient

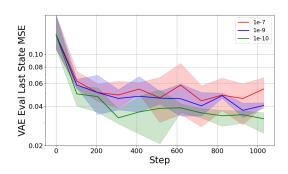


Figure 12: Effect of KL coefficient

A key hyperparameter in LWD is the KL regularization term, $\beta_{\rm KL}$, used during VAE training. In this section, we analyze its impact on the learned latent space using the PushT task with an action horizon of 32. We train three VAEs with $\beta_{\rm KL}$ values of 1e-7, 1e-9, and 1e-10. For evaluation, we sample a trajectory of length 32, encode and decode it via the VAE to generate a policy, and then execute this policy in the environment starting from the same initial state. We compute the MSE between the final state reached after 32 steps and the corresponding state in the original trajectory. Figure 12 in Section D.1 shows this metric across 3 seeds during training. Lower

 β_{KL} values result in lower final-state MSE, indicating better trajectory reconstruction. This is due to a more expressive, multi-modal latent space made possible by weaker regularization, without compromising sampling, as diffusion still operates effectively within this space. Visualizations are provided below in Section D.1. Based on these results, we use $\beta_{KL}=1e-10$ in all PushT experiments.

D.1 KL coefficient ablation latent space

Following the KL ablation experiment above, we analyzed the latent space of the encoded trajectories with PCA, similar to that performed in Section C. The three plots in Figure 13, show that the trajectory encodings get closer and lose behavioral diversity when the KL coefficient is high.

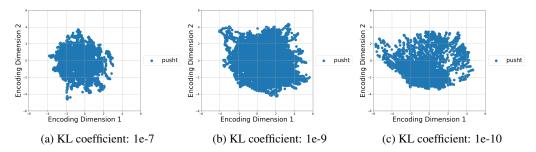


Figure 13: Latent space representation of PushT trajectories at different KL coefficients

E Appendix – Ablation of the Diffusion Model Architecture

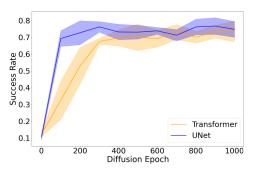


Figure 14: Diffusion Architecture Ablation

Diffusion models typically adopt either UNet-based [26] or Transformer-based [42] architectures (described as medium "m" in Section J.1). To guide our choice for the LWD diffusion policy, we performed an ablation study on the PushT task [9] using an action horizon of 32. As shown in Figure 14, the UNet model demonstrated faster initial learning, achieving higher average success rates early in training. However, both architectures eventually converged to comparable final success rates. For consistency, we adopt the UNet architecture for all other experiments.

F Appendix – Ablation of the Decoder size

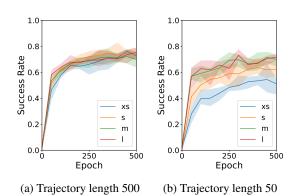


Figure 15: **Effect of VAE decoder size**: For long trajectories, even the smallest decoder (xs) yields high task performance, whereas short trajectories benefit from a larger decoder.

An interesting experiment was the effect of breaking a large trajectory into sub-trajectories and how this affects the latent space. A key takeaway from that experiment was that for halfcheetah locomotion, even small VAE decoders generated accurate policies from trajectory snippets. Whereas, for manipulation tasks from Metaworld, the same-sized small decoder was not capable of reconstructing the original policy. See Section C for this experiment. This finding prompted an ablation on the decoder size, evaluating the average success rate of decoded policies across all 10 Metaworld tasks. Figure 15 illustrates the performance of decoders with varying sizes, denoted as xs (3.9M parameters), s (7.8M parameters), m (15.6M parameters), and l(31.2M parameters). It's important to note

that despite the substantial parameter count of the hypernetwork decoder, the resulting inferred policy remains relatively small (< 100K parameters, see Figure 6). The results demonstrate that increasing the decoder size consistently improves the average success rate of the decoded policies. Refer Section J.3 for more details regarding the decoder size characterization.

This contrasts with rollouts from the HalfCheetah environment, where even smaller decoders generated accurate policies from trajectory snippets. We hypothesize this discrepancy stems from two key factors. First, the cyclic nature of HalfCheetah provides sufficient information within snippets to infer the underlying policy. Second, the increased complexity of Metaworld tasks means that snippets may

lack crucial information for inference. For instance, in a pick-and-place task, a snippet might only capture the "pick" action, leaving the latent without sufficient information to infer the "place" action.

G Appendix – Behavior Reconstruction Analysis

Here, we ask – Does LWD reconstruct the original policies and reproduce diverse behaviors?

G.1 Operator Behavior Analysis

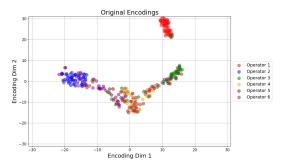


Figure 16: **Behavior distribution** for the Robomimic Lift task

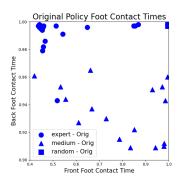
Since LWD assumes trajectory data obtained from a distribution of policies (described in Section 3.1), it offers us access to this distribution through its latent space. As an experiment on the Robomimic Lift task, we encoded trajectories of the demonstrations on this dataset. We set the action horizon to the maximum demonstration length, effectively encoding the entire trajectory for each demonstration. The MH (Multi-Human) dataset contains 50 trajectories from 6 operators, totalling 300 successful trajectories. The operators were varied in proficiency – there were 2 "worse" operators, 2 "okay" operators, and 2 "better" operators. Figure 16 visualizes a

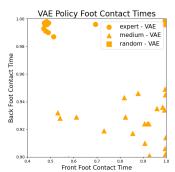
2D t-SNE plot of the latents sampled from the VAE in this experiment. Interestingly, we start to see clusters in this space, each corresponding to a different operator. It is important to note that LWD was never given explicit information regarding these operators, yet was able to cluster them based on their behavior distribution. We believe this interesting result can enable us to filter out unwanted behaviors in imitation learning.

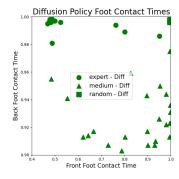
G.2 Locomotion

First, we analyze the behavior reconstruction capability of different components of LWD in locomotion domains. For this experiment, we use the halfcheetah dataset from D4RL [14]. The parameters used for this experiment are shown in Section J.5. Each trajectory in this dataset has a length of 1000. We combine trajectory data from three original behavior policies provided in this dataset: expert, medium, and random. Following [2], we track the foot contact timings of each trajectory as a metric for measuring behavior. For each behavior policy, we get 32 trajectories. These timings are normalized to the trajectory length and are shown in Figure 17. For each plot, the x-axis denotes the foot contact percentage of the front foot, while the y-axis denotes the foot contact percentage of the back foot.

We first visualize the foot contact timings of the original policies in Figure 17a. We see that different running behaviors of the half cheetah can be differentiated in this plot. Then, we train the VAE model on this dataset to embed our trajectories into 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 foot contact timings of these generated policies in Figure 17b. We see that the VAE captures each of the original policy's foot contact distributions, therefore empirically showing that the assumption $p_{\phi_{dec}}(\theta \mid z) = \delta(\theta - f_{\phi_{dec}}(z))$ is reasonable. Then, we train a latent diffusion model conditioned on a behavior specifier (i.e., one task ID per behavior). In Figure 17c, 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. Apart from visual inspection, we also track rewards obtained by the generated policies and empirically calculated Jensen Shannon Divergence between the original and obtained foot contact distributions and observe that LWD maintains behavioral diversity in this locomotion task. See Section H for more details.







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

Figure 17: Foot-contact times shown for various trajectories on the Half Cheetah task. We use foot contact times as the chosen metric to show different behaviors for the half cheetah run task by different policies. The first plot on the left shows the distribution of foot contact percentages for each of the three original policies. The second plot in the center denotes the foot contact percentages for the policies generated by the trained VAE when provided each original policy's entire trajectory. The third plot on the right denotes the foot contact percentages for the policies generated by the diffusion model, trained without any task conditioning.

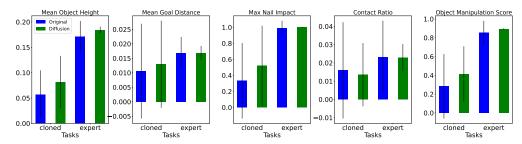


Figure 18: **Behavior Reconstruction for Manipulation**: We track these metrics on the Adroit hammer task, and the LWD-generated policy behaves similarly to the original policy. The 'cloned' bars represent metrics with respect to a human demonstration behavior cloned policy, and 'expert' bars represent metrics from an RL-trained policy.

G.3 Manipulation

To verify the behavior reconstruction capabilities of LWD in manipulation, we also experiment on the D4RL Adroit dataset [45]. We choose a tool use task, where the agent must hammer a nail into a board. We utilize their 5000 expert and 5000 human-cloned trajectories, to train our LWD model. The implementation details are in Section J.6. Then, we evaluate the behavior of the original and generated policy on the following metrics: **Mean object height** - Average height of the object during eval; **Alignment error (goal distance)** - Mean distance between the target and the final goal position; **Max nail impact** - Maximum value of the nail impact sensor during eval; **Contact ratio** - Fraction of time steps where the nail impact sensor value exceeds 0.8; **Object manipulation score** - Proportion of time steps where the object height exceeds 0.04 meters. From Figure 18, we can see that the policy generated by LWD behaves similarly to the original policy.

H Appendix – Behavior Reconstruction Metrics for HalfCheetah

We can analyze the behavior reconstruction capability of LWD by comparing the rewards obtained during a rollout. The VAE parameters used for this experiment are shown in Section J.5. Figure 19 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 2 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 baseline for this experiment, we train a large (5-layer, 512 neurons each) behavior-conditioned MLP on the same mixed dataset with MSE loss. We see that policies generated by LWD consistently achieve a lower JS divergence score than the MLP baseline for expert and medium behaviors. The random behavior is difficult to capture as the actions are almost Gaussian noise. Surprisingly, for the HalfCheetah environment, policies generated by LWD for expert and medium had lower scores than SOS, making it behaviorally indistinguishable from the original policy.

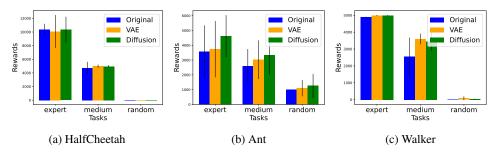


Figure 19: **Reconstruction Rewards**: For each of the 3 environments shown above, the generated policy from trajectory decoded VAE and task-conditioned diffusion model, achieves similar total objective as the original policies. Each bar indicates the mean total objective obtained with error lines denoting the standard deviation.

Environment	Source Policy	Target Policy		
		SOS	MLP	LWD
	Expert	0.187 ± 0.142	1.272 ± 0.911	0.510 ± 0.159
Ant	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
	Expert	0.158 ± 0.146	2.810 ± 1.139	0.088 ± 0.050
HalfCheetah	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
	Expert	0.342 ± 0.329	2.879 ± 1.493	1.093 ± 0.310
Walker2D	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

Table 2: **Behavior Reconstruction**: JS divergence between foot contact distributions from source and target policies. The lower the value, the better.

I Appendix – PushT Execution time comparison

To show the benefit of using LWD over DP when it comes to inference times, we conduct the following experiment. We utilize the models trained in Section 4.1.1 for the PushT task and record the average time taken to complete 256 environment steps at different action horizons. Figure 20 plots the minimum eval time that can achieve the corresponding success rate in the x-axis. This time was recorded while running on an NVIDIA 3090 GPU. The lower the area under this plot, the better. We record these metrics for varying perturbation levels. We see that for all settings, we achieve corresponding success rates with lower evaluation time. Thus showing the benefit of longer action horizons with LWD.

Crucially, we observed that episodes taking 1 second for a rollout had only 2 diffusion queries (large action horizon), whereas those taking 6 seconds had 32 diffusion queries (small action horizon). The number of environment steps is the same. Thus, it is clear that diffusion-based methods can

consume up most of the total rollout time, significantly exceeding the time spent interacting with the environment. This substantial computational overhead underscores the critical need for optimizing the inference efficiency of diffusion models, particularly in time-critical robotics applications.

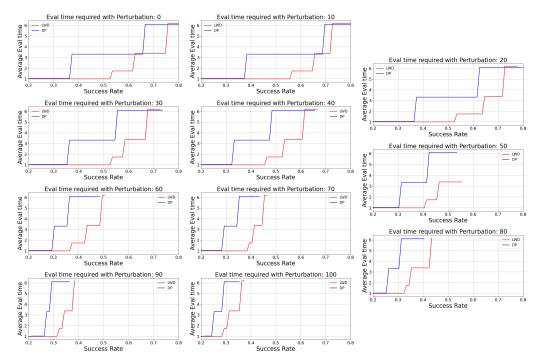


Figure 20: Average minimum episode eval time required to achieve corresponding success rate.

J Appendix – Implementation Details

The following are the hyperparameters we use for our experiments:

J.1 Baseline Diffusion Policy model

To train the diffusion policy baseline model shown in Figure 6, we utilize the training script provided by the authors of DP here:

https://colab.research.google.com/drive/1gxdkgRVfM55zihY9TFLja97cSVZOZq2B?usp=sharing. To set the model size we use the following parameters:

Size	Diffusion Step Embed Dim	Down Dims	Kernel Size
extra-small: (s)	64	[16, 32, 64]	5
small: (s)	256	[32, 64, 128]	5
large: (m)	256	[128, 256, 256]	5
large: (1)	256	[256, 512, 1024]	5
extra large: (xl)	512	[512, 1024, 2048]	5

Table 3: Architectural configurations for the ConditionUnet1D Diffusion Policy (DP) across different model sizes.

For the ablation described in Section E, we use a transformer architecture, the details of which are:

J.2 VAE Encoder details

For the encoder, we first flatten the trajectory to form a one-dimensional array, which is then fed to a Multi-Layer Perceptron with three hidden layers of 512 neurons each.

Size	Diffusion Step Embed Dim	Model Dim	# Layers	# Heads
extra-small: (xs)	64	64	3	2
small: (s)	128	128	4	4
medium: (m)	256	256	6	8
large: (1)	256	512	8	8
extra-large: (xl)	512	768	12	12

Table 4: Architectural configurations for Transformer-based Diffusion models across different model sizes.

J.3 VAE Hypernetwork decoder size characterization

For the hypernetwork, we utilize an HMLP model (a full hypernetwork) from the https://hypnettorch.readthedocs.io/en/latest/ package with default parameters. We condition the HMLP model on the generated latent of dimension 256. To vary the size of the decoder, as explained in Section F, we set the hyperparameter in the HMLP as shown in Table 5

Size	No. of parameters	layers
XS	3.9M	[50, 50]
S	7.8M	[100, 100]
m	15.6 M	[200, 200]
1	31.2M	[400, 400]

Table 5: VAE size varying parameters

J.4 Diffusion model parameters

For all our experiments, we utilize the same ConditionalUnet1D network from [9] as the diffusion model. This is the same as the DP-medium (m) model described in Section J.1.

J.5 Mujoco locomotion tasks

We use the following hyperparameters to train VAEs for all D4RL mujoco tasks shown in the paper. To show the effect of shorter trajectories in Section C, we change the Trajectory Length to 100.

Parameter	Value
Trajectory Length	1000
Batch Size	32
VAE Num Epochs	150
VAE Latent Dimension	256
VAE Decoder Size	s
Evaluation MLP Layers	{256, 256}
VAE Learning Rate	3×10^{-4}
KL Coefficient	1×10^{-6}
Diffusion Num Epochs	200

Table 6: Mujoco locomotion hyperparameters.

J.6 Adroit Hammer task

We use the same hyperparameters as Table 6 and override the following hyperparameters to train VAEs for the D4RL Adroit hammer task shown in the paper.

Further, for the experiment where we show the hammer task can be composed of sub-tasks, we change the Trajectory Length to 32 to enable LWD to learn the distribution of shorter horizon policies.

Parameter	Value
Trajectory Length	128
VAE Num Epochs	20
Diffusion Num Epochs	10

Table 7: Adroit hammer hyperparameters.

J.7 PushT and Robomimic LWD

For all the experiments shown in Section 4.1.1, we use the same hyper-parameters described in Table 6, and override the following:

Parameter	Value
Trajectory Length	256
VAE Num Epochs	1000
Diffusion Num Epochs	1000
Diffusion Model size	1
VAE Decoder Size	1
VAE KL coefficient	1e - 10

Table 8: PushT LWD hyperparameters.

J.8 Metaworld tasks

For all the experiments shown in Section 4.1.2, we use the same hyper-parameters described in Table 6, and override the following:

Parameter	Value
Trajectory Length	500
VAE Num Epochs	100
Diffusion Num Epochs	100
VAE Decoder Size	xs

Table 9: Metaworld hyperparameters.

To show the effect of shorter trajectories in Section C, we change the Trajectory Length to 50.

J.9 Compute Resources

Each VAE and diffusion experiment was run on jobs that were allocated 6 cores of a Intel(R) Xeon(R) Gold 6154 3.00GHz CPU, an NVIDIA GeForce RTX 2080 Ti GPU, and 108 GB of RAM.

K Societal Impact

This paper presents work intending to advance the field of Machine Learning in Robotics. The authors are committed to ensuring that the research is used for the benefit of society and does not contribute to any harm or negative consequences. While there are many potential societal consequences robotics and machine learning in general, there are none that are specific to this work alone, that we can foresee.