
STITCH-OPE: Trajectory Stitching with Guided Diffusion for Off-Policy Evaluation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Off-policy evaluation (OPE) estimates the performance of a target policy using
2 offline data collected from a behavior policy, and is crucial in domains such as
3 robotics or healthcare where direct interaction with the environment is costly or
4 unsafe. Existing OPE methods are ineffective for high-dimensional, long-horizon
5 problems, due to exponential blow-ups in variance from importance weighting or
6 compounding errors from learned dynamics models. To address these challenges,
7 we propose STITCH-OPE, a model-based generative framework that leverages
8 denoising diffusion for long-horizon OPE in high-dimensional state and action
9 spaces. Starting with a diffusion model pre-trained on the behavior data, STITCH-
10 OPE generates synthetic trajectories from the target policy by guiding the denoising
11 process using the score function of the target policy. STITCH-OPE proposes
12 two technical innovations that make it advantageous for OPE: (1) prevents over-
13 regularization by subtracting the score of the behavior policy during guidance,
14 and (2) generates long-horizon trajectories by stitching partial trajectories together
15 end-to-end. We provide a theoretical guarantee that under mild assumptions, these
16 modifications result in an exponential reduction in variance versus long-horizon
17 trajectory diffusion. Experiments on the D4RL and OpenAI Gym benchmarks show
18 substantial improvement in mean squared error, correlation, and regret metrics
19 compared to state-of-the-art OPE methods.

20 1 Introduction

21 Given the slow and risky nature of online data collection, real-world applications of reinforcement
22 learning often require offline data for policy learning and evaluation [22, 38]. An important problem
23 of working with offline data is *off-policy evaluation* (OPE), which aims to evaluate the performance
24 of target policies using offline data collected from other behavior policies. One practical advantage of
25 OPE is that it saves the cost of evaluation on hardware in embodied applications in the real world [27].
26 However, a central challenge of OPE is the presence of *distribution shift* induced by differences in
27 behavior and target policies [22, 3]. This can lead to inaccurate estimates of policy values, making it
28 difficult to trust or select between multiple target policies before they are deployed [43, 25].

29 Numerous approaches have attempted to address the distribution shift in offline policy evaluation
30 by reducing either the variance of the policy value or its bias, but they are typically ineffective in
31 high-dimensional long-horizon problems. For example, Importance Sampling (IS) [32] estimates
32 the value of the target policy by weighing the behavior policy rollouts according to the ratio of their
33 likelihoods. However, it suffers from the so-called *curse of horizon* where the variance of the estimate
34 increases exponentially in the evaluation horizon [24]. More recent model-free OPE estimators
35 reduce or eliminate the explosion in variance by estimating the long-run state-action density ratio
36 $d^\pi(s, a)/d^\beta(s, a)$ between the target and behavior policy [24, 30, 42], yet they have demonstrated

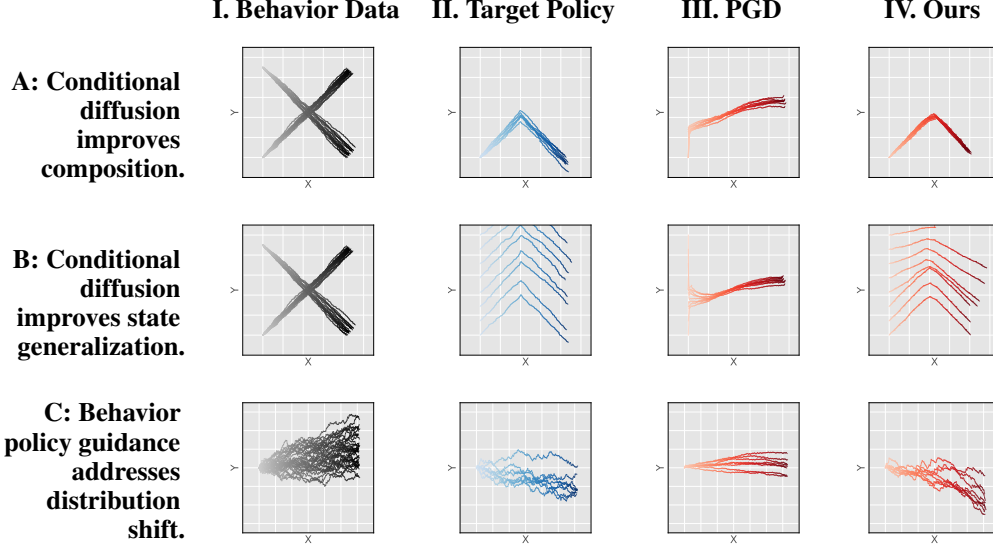


Table 1: A 2D toy problem with Gaussian dynamics illustrates the advantages of STITCH-OPE. **Row A:** Behavior data is a mixture of two datasets generated by different behavior policies β_1 and β_2 . The target policy in column II is a piecewise function following β_1 on the left half-space and β_2 on the right half-space. Policy-Guided Diffusion (PGD) in column III estimates the behavior trajectory distribution by training a diffusion model on the behavior dataset, then leverages guided diffusion to generate target policy trajectories [15]. PGD is unable to stitch behavior trajectories correctly while STITCH-OPE is. One explanation is that conditional diffusion provides a higher entropy sampling distribution than full-length diffusion, and thus ensures a broader coverage of the modes in the behavior data (see Section 3.3 for details). **Row B:** The initial state is varied during trajectory generation. As shown in column III, PGD cannot generalize to new initial states. On the contrary, STITCH-OPE is trained on sub-trajectories that start in arbitrary states in the behavior dataset as opposed to only initial states, which allows it to achieve better generalization. **Row C:** A scenario with severe distribution shift is presented, where the behavior and target policies move the agent in different directions. As shown in columns III and IV, the negative behavior guidance term is essential to prevent over-regularization, which can prevent guided diffusion from addressing distribution shift and lead to biased value estimation.

poor empirical performance on high-dimensional tasks where the behavior and target policies are different (i.e. the behavior policy is not a noisy version of the target policy) [13].

As an alternative approach, model-based OPE estimators typically learn an empirical autoregressive model of the environment and reward function from the behavior data, which is used to generate synthetic rollouts from the target policy for offline evaluation [18, 37, 45]. Some advantages of the model-based paradigm include sample efficiency [23], exploitation of prior knowledge about the dynamics [11], and better generalization to unseen states [44]. Although model-based OPE methods often scale well to high-dimensional short-horizon problems – owing to the scalability of the deep model-based RL paradigm – their robustness diminishes in long-horizon tasks due to the compounding of errors in the approximated dynamics model [13, 16, 17].

Driven by the recent successes of generative diffusion in RL [28, 47, 1, 15, 29, 34], we propose *Sub-Trajectory Importance-Weighted Trajectory Composition for Long-Horizon OPE* for model-based off-policy evaluation in long-horizon high-dimensional problems. STITCH-OPE first trains a diffusion model on behavior data, allowing it to generate dynamically feasible behavior trajectories [17]. STITCH-OPE differs from prior work by training the diffusion model on short sub-trajectories instead of full rollouts, where sub-trajectory generation is conditioned on the final state of the previous generated sub-trajectory. This enables accurate trajectory “stitching” using short-horizon rollouts, while minimizing compounding error of full-trajectory rollouts, thus bridging the gap between model-based OPE and full-trajectory offline diffusion.

STITCH-OPE explicitly accounts for distribution shift in OPE by guiding the diffusion denoising process [9, 17] during inference. This can be achieved by selecting the guidance function to be the difference between the score functions of the target and behavior policies. A significant advantage of guided diffusion is that it eliminates the need to retrain the diffusion model for each new target policy. By pretraining the model on a variety of behavior datasets, generalization can be achieved during guided sampling to produce feasible trajectories under the target policy, leading to robust off-policy estimates for target policies that lack offline data. STITCH-OPE contributes the following novel technical innovations that we consider critical to the successful and robust application of diffusion models for OPE:

- **Trajectory Stitching with Conditional Diffusion.** We propose a state-conditioned guided diffusion model for generating short sub-trajectories from the target policy (Figure 1). This significantly improves the quality (Table 1, row A) and generalization (Table 1, row B) of trajectory generation in long-horizon tasks. Theorem 3.3 also provides theoretical bounds on the bias and variance of our proposed approach.
- **Behavior Guidance.** We show that the negative score function of the behavior policy mitigates the diffusion model from collapsing to trajectories with large behavior likelihood, thus improving generalization out-of-data (Table 1, row C). This negative score function naturally arises from viewing the likelihood ratio (i.e. in importance sampling) as the classification density in diffusion guidance [9] – a connection missed in prior work [15].
- **Robustness Across Problem Difficulty.** Finally, we evaluate STITCH-OPE on the OpenAI Gym control suite [4] and the D4RL offline RL suite [12], showing significant improvements compared to other recent OPE estimators across a variety of metrics (mean squared error, rank correlation and regret), problem dimension and evaluation horizon. To our knowledge, STITCH-OPE is the first work to demonstrate robust off-policy evaluation in high-dimensional long-horizon tasks.

2 Preliminaries

Markov Decision Processes. We consider the standard *Markov Decision Process* (MDP), which consists of a 6-tuple $\langle \mathcal{S}, \mathcal{A}, R, P, d_0, \gamma \rangle$ [33] where: \mathcal{S} is the continuous state space, \mathcal{A} is the continuous action space, R is the reward function, P is the Markov transition probability distribution, d_0 is the initial state distribution and $\gamma \in [0, 1]$ is the discount factor.

A policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, \infty)$ observes the current state s , samples an action according to its conditional distribution $\pi(\cdot|s)$, and observes the immediate reward $R(s, a)$ and the next state $s' \sim P(\cdot|s, a)$. The interaction of a policy with an MDP generates a set of trajectories of states, actions, and rewards. The goal of reinforcement learning is to learn a policy π that maximizes the expected return:

$$J(\pi) = \mathbb{E}_{\tau \sim p_\pi} \left[\sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right], \quad p_\pi(\tau) = d_0(s_0) \prod_{t=0}^{T-1} \pi(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

over length- T trajectories $\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$ induced by policy π .

Off-Policy Evaluation. The goal of *Off-Policy Evaluation* (OPE) is to estimate the expected return of some *target policy* π given only a data set of trajectories \mathcal{D}_β from some *behavior policy* β . To estimate $J(\pi)$, it is necessary to approximate the distribution over trajectories $p_\pi(\tau)$ induced by π using only samples from the distribution $p_\beta(\tau)$. Hence, the phenomenon of *distribution shift* arises whenever β is sufficiently different from π , in which case $p_\beta(\tau)$ is not a suitable proxy for obtaining samples from $p_\pi(\tau)$ and corrections must be made to account for the distribution shift.

Denoising Diffusion Models. *Denoising Diffusion Probabilistic Models* (DDPMs) [35, 14] are a class of generative models to sample from a given distribution. Given a dataset $\{x_i\}_{i=1}^N$ and the K -step (forward) noise process $x_i^k = \sqrt{\alpha_k} x_i^{k-1} + \sqrt{1 - \alpha_k} \epsilon$, where $\epsilon \sim \mathcal{N}(0, I)$ is independent, DDPMs are trained to perform the K -step (backward) denoising process to recover x_i from $x_i^K \sim \mathcal{N}(0, I)$. This is accomplished by running the forward process $x_i^0 \rightarrow x_i^k$ on the original x_i and training the DDPM ϵ_θ on the denoising process $x_i^{k-1} | x_i^k \sim \mathcal{N}(\mu(x_i^k, k), \sigma_k^2 I)$ to predict the noise ϵ so that the distribution over denoised samples x_i^0 matches x_i . The standard reparameterization

105 $\mu(x_i^k, k) = (x_i^k - \epsilon_\theta(x_i^k, k)(1 - \alpha_k)/\sigma_k)/\sqrt{\alpha_k}$ maps x^k and the predicted noise $\epsilon_\theta(x_i^k, k)$ to x^{k-1}
 106 to undo the noise accumulated during the forward process. The following loss function is typically
 107 used to train a diffusion model [14]

$$\mathcal{L}(\theta) = \mathbb{E}_{k, x_i, \epsilon \sim \mathcal{N}(0, I)} [\|\epsilon - \epsilon_\theta(x_i^k, k)\|^2]. \quad (1)$$

108 **Guided Diffusion.** It is possible to guide the sampling from a trained diffusion model to maximize
 109 some classifier $p(y|x)$ [9]. A key observation of diffusion is that the backward diffusion process can
 110 be well approximated by a Gaussian when the noise is small, that is, $x^k|x^{k+1} \sim \mathcal{N}(\mu_k, \Sigma_k)$. Next,
 111 observe that $p(x^k|x^{k+1}, y) \propto p(x^k|x^{k+1})p(y|x^k)$ by Bayes' rule. Applying the first-order Taylor
 112 approximation $\log p(y|x^k) \approx \log p(y|\mu_k) + (x^k - \mu_k)g(\mu_k)$ at μ_k , where $g(u) = \nabla_x \log p(y|x)|_{x=u}$,
 113 it can be shown that:

$$\begin{aligned} \log(p(x^k|x^{k+1})p(y|x^k)) &\propto -\frac{1}{2}(x^k - \mu_k - \Sigma_k g(\mu_k))^T \Sigma_k^{-1} (x^k - \mu_k - \Sigma_k g(\mu_k)) \\ &\propto \log \mathcal{N}(x^k; \mu_k + \Sigma_k g(\mu_k), \Sigma_k). \end{aligned} \quad (2)$$

114 In other words, it is possible to sample from the conditional (guided) distribution $p(x^k|x^{k+1}, y)$ by
 115 sampling from the original diffusion model with its mean shifted to $\mu_k + \Sigma_k g(\mu_k)$. Appendix A
 116 provides a worked example illustrating guided diffusion for a mixture of Gaussians.

117 3 Proposed Methodology

118 The *direct method* for off-policy evaluation [10] estimates the single-step autoregressive model
 119 $\hat{P}(s_t|s_{t-1}, a_{t-1})$ and the reward function $\hat{R}(s_t, a_t)$ from the behavior data. Then, it draws tar-
 120 get policy trajectories $\tau \sim p_\pi(\tau)$ by forward sampling, that is, $s_0 \sim d_0, a_0 \sim \pi(\cdot|s_0), s_1 \sim$
 121 $\hat{P}(\cdot|s_0, a_0), \dots, s_T \sim \hat{P}(\cdot|s_{T-1}, a_{T-1})$. However, even small errors in \hat{P} can lead to significant bias
 122 in $J(\pi)$ due to the compounding of errors over long horizon T [19, 16]. STITCH-OPE avoids the
 123 compounding problem by generating the partial trajectory in a single (backward diffusion) pass,
 124 leading to more accurate OPE estimates over a long horizon.

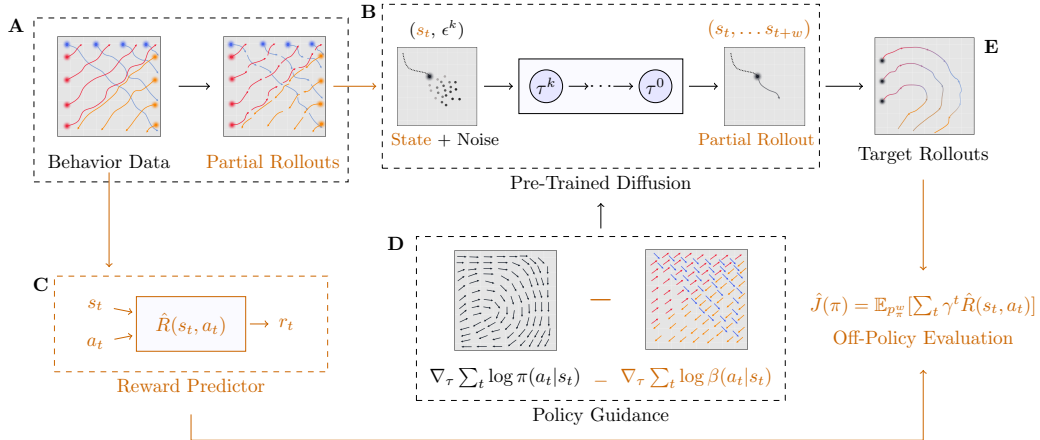


Figure 1: A conceptual illustration of STITCH-OPE, with novel contributions highlighted in orange. **A:** Behavior data is sliced into partial trajectories of length w . **B:** The data is fed to a conditional diffusion model taking a w -length sequence of Gaussian noise ϵ and state s_t as inputs, and applies the backward diffusion process to predict the behavior trajectory of length w beginning in state s_t . **C:** To evaluate policies, STITCH-OPE also trains a neural network on the behavior transitions to predict the immediate reward. **D:** It then applies guided diffusion on the pretrained diffusion model to generate a batch of partial target trajectories of length w , where the guidance function incorporates the score function of the target policy and the behavior policy. **E:** The guided partial trajectories are stitched end-to-end to produce full-length target trajectories. Finally, the guided trajectories are evaluated using the empirical reward function $\hat{R}(s, a)$, and averaged to estimate the value of the target policy.

3.1 Guided Diffusion for Off-Policy Evaluation

It is possible to approximate p_π using guided diffusion by interpreting each data point x_i as a full trajectory τ . Given a behavior policy β and corresponding length- T trajectory distribution $p_\beta(\tau)$, the corresponding length- T trajectory distribution of target policy π can be written as:

$$p_\pi(\tau) = d_0(s_0) \prod_{t=0}^{T-1} \beta(a_t|s_t) P(s_{t+1}|s_t, a_t) \frac{\pi(a_t|s_t)}{\beta(a_t|s_t)} = p_\beta(\tau) \prod_{t=0}^{T-1} \frac{\pi(a_t|s_t)}{\beta(a_t|s_t)}, \quad (3)$$

which is the standard importance sampling correction [32]. We address the question of tractably learning $p_\beta(\tau)$ by training a diffusion model $\hat{p}_\beta(\tau)$ on the offline behavior data set \mathcal{D}_β [17], thus approximating $\hat{p}_\beta(\tau) \approx p_\beta(\tau)$. Specifically, the diffusion model learns to map a trajectory consisting of pure noise, $\tau^k = (s_0^k, a_0^k, \dots, s_T^k)$, to a noiseless behavior trajectory $\tau^0 = (s_0^0, a_0^0, \dots, s_T^0)$.

A key observation is that we can bypass importance sampling in (3) by guiding the generation process $\hat{p}_\beta(\tau)$ towards $p_\pi(\tau)$ using diffusion guidance (2) [17]. Specifically, let $x^k = \tau^k$ denote a noisy behavior trajectory at step k of the forward diffusion process, and let $y \in \{0, 1\}$ be a binary outcome with $p(y = 1|\tau) \propto \prod_{t=0}^{T-1} \frac{\pi(a_t|s_t)}{\beta(a_t|s_t)}$. Intuitively, y indicates whether the trajectory τ is generated by the target policy π ($y = 1$) or the behavior policy β ($y = 0$), and the likelihood ratio determines the odds that $y = 1$ given τ . By (3),

$$p_\pi(\tau) \propto p_\beta(\tau) p(y = 1|\tau),$$

and thus the backward diffusion process for generating target policy trajectories for OPE can be approximated with guidance (2):

$$\begin{aligned} \log p_\pi(\tau^k|\tau^{k+1}) &\propto \log(p_\beta(\tau^k|\tau^{k+1}) p(y = 1|\tau^{k+1})) \\ &\approx \log \mathcal{N}(\tau^{k+1}; \mu_k + \Sigma_k \nabla_\tau \log p(y = 1|\tau)|_{\tau^{k+1}}, \Sigma_k), \end{aligned} \quad (4)$$

where $p_\beta(\tau^k|\tau^{k+1}) = \mathcal{N}(\mu_k, \Sigma_k)$ is the backward diffusion process. Therefore, we can obtain feasible target policy trajectories using the guidance function:

$$g(\tau) = \nabla_\tau \log p(y = 1|\tau) = \nabla_\tau \sum_{t=0}^{T-1} \log \pi(a_t|s_t) - \nabla_\tau \sum_{t=0}^{T-1} \log \beta(a_t|s_t). \quad (5)$$

Given the approximate sampling distribution over the trajectories of the target policy described above, $\hat{p}_\pi(\tau^0) = \int \dots \int \mathcal{N}(\tau^K; 0, I) \prod_{k=1}^K p_\pi(\tau^{k-1}|\tau^k) d\tau^K \dots d\tau^1$, and an empirical reward function $\hat{R}(s, a)$, it is straightforward to estimate the expected return (or a statistic such as variance or quantile) given any target policy, i.e. $\hat{J}(\pi) = \mathbb{E}_{\tau \sim \hat{p}_\pi} \left[\sum_t \gamma^t \hat{R}(s_t, a_t) \right] \approx J(\pi)$.

3.2 Negative Behavior Guidance

The target policy score function, $g_{\text{simple}}(\tau) = \nabla_\tau \sum_{t=0}^{T-1} \log \pi(a_t|s_t)$ [15], provides a simple guidance function for OPE. However, it corresponds to a biased estimator of $p_\pi(\tau)$ in the context of (3) and can generate trajectories that are unlikely under the target policy, as illustrated using the GaussianWorld domain in Table 1 (see Appendix B for details). The behavior policy β returns a positive angle in each state and the target policy π returns a negative angle, to test the performance of both guidance functions under distribution shift. Conclusions are summarized in row C of Table 1. The omission of the negative guidance term results in a sampling distribution that collapses to a high-density region under $p_\beta(\tau)$, where $p_\pi(\tau)$ could be small. In other words, **behavior guidance prevents the guided sampling distribution \hat{p}_π from becoming over-regularized.**

In our empirical evaluation, we employ the following generalization of (5) to allow fine-grained control over the relative importance of the target and behavior policy guidance

$$g(\tau) = \alpha \nabla_\tau \sum_{t=0}^{T-1} \log \pi(a_t|s_t) - \lambda \nabla_\tau \sum_{t=0}^{T-1} \log \beta(a_t|s_t). \quad (6)$$

Ignoring the normalizing constant which does not dependent on τ , (6) is equivalent to sampling from the following re-weighted trajectory distribution

$$q_\pi(\tau) \propto p_\beta(\tau) \prod_{t=0}^{T-1} \frac{\pi(a_t|s_t)^\alpha}{\beta(a_t|s_t)^\lambda}, \quad (7)$$

which can be interpreted as a *tempered posterior distribution* [2, 5] over trajectories; the importance of the likelihood terms associated with π and β are controlled by the choice of α and λ , respectively. Note that the choice $\alpha = \lambda = 1$ reduces to the standard guidance function $g(\tau)$. This is not a good empirical choice because it can push the backward diffusion process too far from the behavior distribution, leading to infeasible trajectories or instability. Instead, typical choices satisfy $\lambda < \alpha$ (see Appendix L.1 for an additional experiment confirming this). The choice $\alpha = 1, \lambda = 0$ reduces to $g_{simple}(\tau)$ and is unsuitable for OPE.

3.3 Sub-Trajectory Stitching with Conditional Diffusion

Recent work has shown that full-length diffusion models do not provide sufficient compositionality for accurate long-horizon sequence generation [6]. In addition, full-length prediction requires the generation of sequences of length $T \cdot (\dim(\mathcal{A}) + \dim(\mathcal{S}))$; this may be infeasible or inefficient on resource-constrained systems, when T is large or when \mathcal{A} or \mathcal{S} is high-dimensional.

To tackle these limitations, STITCH-OPE trains a conditional diffusion model to generate behavior sub-trajectories of length $w \ll T$. To allow for a more flexible composition of behavior trajectories during guidance, generation in STITCH-OPE is performed in a semi-autoregressive manner from the diffusion model, which is conditioned on the last state of the previously generated sub-trajectory.

Specifically, the conditional diffusion model in STITCH-OPE, denoted as $\epsilon_\theta(\tau_{t:t+w}^k, k | s_t^0)$, denoises a length- w noisy sub-trajectory $\tau_{t:t+w}^k = (s_t^k, a_t^k, s_{t+1}^k, a_{t+1}^k, \dots, s_{t+w-1}^k, a_{t+w-1}^k, s_{t+w}^k)$ conditioned on the last state s_t^0 of the previously generated sub-trajectory $\tau_{t-w:t}^0$. Generalizing (1), the loss function of STITCH-OPE is thus

$$\mathcal{L}_{STITCH-OPE}(\theta) = \mathbb{E}_{k, t, \tau_{t:t+w} \sim \mathcal{D}_\beta, \epsilon \sim \mathcal{N}(0, I)} [\|\epsilon - \epsilon_\theta(\tau_{t:t+w}^k, k | s_t^0)\|^2].$$

Next, writing $p_\beta(\tau_{t:t+w} | s_t^0)$ to denote the sampling distribution over fully denoised sub-trajectories $\tau_{t:t+w}^0$ conditioned on s_t^0 , the sampling process (3) of STITCH-OPE can be written as:

$$p_\pi^w(\tau) = \prod_{t=0}^{T/w-1} \left(p_\beta(\tau_{wt:w(t+1)} | s_{wt}^0) \cdot \prod_{u=wt}^{w(t+1)-1} \frac{\pi(a_u | s_u)}{\beta(a_u | s_u)} \right), \quad (8)$$

where target trajectories are generated by guiding the conditional diffusion model (analogous to (4)) according to

$$g(\tau_{wt:w(t+1)}) = \alpha \nabla_{\tau_{wt:w(t+1)}} \sum_{u=wt}^{w(t+1)-1} \log \pi(a_u | s_u) - \lambda \nabla_{\tau_{wt:w(t+1)}} \sum_{u=wt}^{w(t+1)-1} \log \beta(a_u | s_u).$$

A complete algorithm description of STITCH-OPE is provided in Appendix E.

To understand the intuition that the conditional diffusion model offers better compositionality than the full-horizon prediction, we decompose the behavior trajectory distribution as a mixture over the trajectories τ_j in \mathcal{D}_β :

$$p_\beta(s_t, a_t, \dots, s_{T-1} | s_0, a_0, \dots, s_t) \approx \sum_{\tau_j \in \mathcal{D}_\beta} p_\beta(s_t, a_t, \dots, s_{T-1} | s_t, \tau_j) p(\tau_j | s_0, a_0, \dots, s_t).$$

Meanwhile, the conditional diffusion model ignores the full history of past states, i.e.:

$$p_\beta(s_t, a_t, \dots, s_{T-1} | s_0, a_0, \dots, s_t) \approx \sum_{\tau_j \in \mathcal{D}_\beta} p_\beta(s_t, a_t, \dots, s_{T-1} | s_t, \tau_j) p(\tau_j | s_t).$$

$p(\tau_j | s_t)$ has higher entropy than $p(\tau_j | s_0, a_0, \dots, s_t)$ since it is conditioned on less information (see Appendix C for a proof), and thus provides a broader coverage of the diverse modes in the behavior dataset. This improves the compositionality of guided long-horizon trajectory generation. Row A of Table 1 illustrates this claim empirically using the GaussianWorld problem. A further claim is that the STITCH-OPE model can generalize better across initial states with low, or even zero, probability under d_0 (see row B of Table 1). This occurs because p_β^w is trained on sub-trajectories starting in arbitrary states in \mathcal{D}_β , as opposed to states sampled only from d_0 . Therefore, **sliding windows strike an optimal balance between autoregressive methods and full-length trajectory diffusion, providing good compositionality while avoiding the error compounding in terms of T .**

3.4 Theoretical Analysis

We provide theoretical guarantees for our proposed STITCH-OPE method by analyzing its bias and variance. The first prerequisite assumption is standard in OPE [41, 26] and limits the ratio π/β .

Assumption 3.1. There is a constant κ such that $\frac{\pi(a|s)}{\beta(a|s)} \leq \kappa$ for all $s \in \mathcal{S}, a \in \mathcal{A}$.

The second prerequisite assumption bounds the total variation between the learned length- w trajectory distribution \hat{p}_β^w and the true distribution p_β^w under the behavior policy β .

Assumption 3.2. $TV(p_\beta^w, \hat{p}_\beta^w) \leq \delta_\beta$ for some constant δ_β .

Our main result is a bound on the mean squared error of the STITCH-OPE estimator. We defer the full proofs, technical lemmas, and definitions to Appendix D.

Theorem 3.3. Define \hat{p}_π as the (length- T) trajectory distribution of the guided diffusion model, and p_π as the true trajectory distribution under the target policy π . Under Assumptions 3.1 and 3.2, the mean squared error (MSE) of the STITCH-OPE return \hat{J} satisfies:

$$\mathbb{E}_{\hat{p}_\pi} \left[(\hat{J} - J(\pi))^2 \right] \leq \underbrace{\left(\frac{2B_w}{1-\gamma^w} \kappa^w \delta_\beta \right)^2}_{\text{Bias}^2} + 10 \underbrace{\left(\frac{T}{w} \right)^2 B_w^2 \kappa^w \delta_\beta + \frac{8B_w^2}{1-\gamma^{2w}} \kappa^w \delta_\beta + \text{Var}_{p_\pi}(J)}_{\text{Variance}},$$

where $B_w = \frac{1-\gamma^w}{1-\gamma} \sup_{s,a} |R(s, a)|$ is a bound on the maximum length- w discounted return, and J is the return under p_π .

Remarks. Theorem 3.3 resembles the $O(\exp(cT))$ bound of IS-based methods [24, 26], but is expressed in terms of w rather than T (with a more favorable $O(T^2)$ dependence). Since w is a fixed hyper-parameter typically chosen to be much smaller than T in practice, **STITCH-OPE provides an exponential reduction in MSE versus both importance sampling and length- T diffusion!** In Section 4, we validate this claim further by showing that STITCH-OPE outperforms full trajectory diffusion (PGD) on most benchmarks (see Appendix L.2 for an additional experiment confirming small $w > 1$ is ideal in practice). The MSE decreases as the error in the learned behavior model δ_β decreases. In practice, δ_β is easier to estimate and control than the error of the target density p_π^w . It is also important to note that the variance cannot be reduced below $\text{Var}_{p_\pi}(J)$, the intrinsic variance of the environment and the target policy.

4 Empirical Evaluation

Our empirical evaluation aims to answer the following research questions:

1. Does the combination of conditional diffusion and negative guidance (as hypothesized in Table 1) translate to robust OPE performance on standard benchmarks?
2. Is STITCH-OPE robust across problem size (e.g., state/action dimension, horizon)?
3. Is STITCH-OPE robust across different levels of optimality of the target policy and the classes of policies?

4.1 Experiment Details

Domains We evaluate the performance of STITCH-OPE in high-dimensional long-horizon tasks using the standard D4RL benchmark [12] and their respective benchmark policies [13]. Specifically, we use the `halfcheetah-medium`, `hopper-medium` and `walker2d-medium` behavior datasets. Each evaluation consists of 10 target policies $\pi_1, \pi_2, \dots, \pi_{10}$ trained at varying levels of ability [13]. We also carry out similar experiments using classical control tasks (Pendulum and Acrobot) from OpenAI Gym [4], to evaluate the competitiveness of STITCH-OPE on standard benchmarks on which other baselines have been extensively evaluated. For this set of environments, we obtain the target policies by running the twin-delayed DDPG algorithm [8] (see Appendix G for details). We set the trajectory length to $T = 768$ for all D4RL problems, $T = 256$ for Acrobot, and $T = 196$ for Pendulum. We also use $\gamma = 0.99$ in all experiments. The domain details are provided in Appendix F, and the training details of STITCH-OPE are provided in Appendix J.

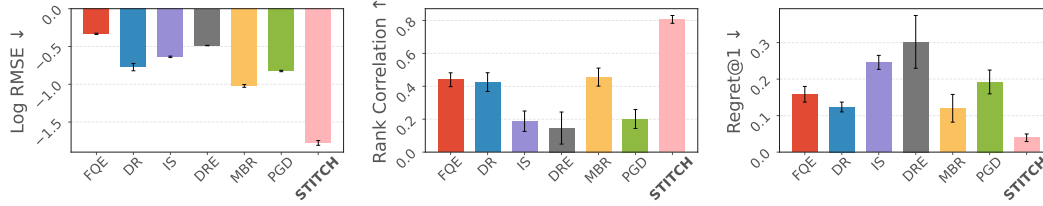


Figure 2: Mean overall performance of all baselines, averaged across environments. Error bars represent \pm one standard error.

		FQE	DR	IS	DRE	MB	PGD	Ours
Log RMSE ↓	Hopper	-0.42 ± 0.03	-0.57 ± 0.02	-0.48 ± 0.01	-0.42 ± 0.00	-1.70 ± 0.04	-1.22 ± 0.02	-2.33 ± 0.02
	Walker2d	-0.48 ± 0.01	-1.25 ± 0.08	-0.71 ± 0.01	-0.45 ± 0.00	-0.88 ± 0.01	-0.32 ± 0.01	-1.33 ± 0.01
	HalfCheetah	-0.05 ± 0.00	0.01 ± 0.01	-0.84 ± 0.02	-1.19 ± 0.00	-0.37 ± 0.00	-1.47 ± 0.00	-0.85 ± 0.01
	Pendulum	-0.58 ± 0.00	-1.02 ± 0.04	-0.15 ± 0.00	-0.58 ± 0.00	-0.43 ± 0.01	-0.91 ± 0.01	-2.34 ± 0.07
	Acrobot	-0.14 ± 0.00	-0.49 ± 0.06	-1.00 ± 0.01	0.20 ± 0.01	-1.54 ± 0.02	-0.13 ± 0.01	-2.02 ± 0.05
Rank Corr. ↑	Hopper	0.17 ± 0.05	0.69 ± 0.06	-0.06 ± 0.13	-0.09 ± 0.14	0.52 ± 0.03	0.36 ± 0.09	0.76 ± 0.02
	Walker2d	0.41 ± 0.05	0.50 ± 0.02	0.51 ± 0.11	0.42 ± 0.05	0.65 ± 0.04	-0.07 ± 0.10	0.63 ± 0.03
	HalfCheetah	-0.03 ± 0.06	-0.48 ± 0.07	0.57 ± 0.06	0.80 ± 0.02	0.32 ± 0.03	0.50 ± 0.00	0.87 ± 0.01
	Pendulum	0.89 ± 0.03	0.72 ± 0.07	-0.60 ± 0.00	-0.40 ± 0.15	0.84 ± 0.06	0.54 ± 0.02	0.96 ± 0.02
	Acrobot	0.75 ± 0.02	0.63 ± 0.08	0.52 ± 0.01	0.01 ± 0.12	0.53 ± 0.11	0.43 ± 0.14	0.82 ± 0.04
Regret@1 ↓	Hopper	0.13 ± 0.03	0.05 ± 0.02	0.13 ± 0.02	0.27 ± 0.17	0.04 ± 0.03	0.04 ± 0.01	0.11 ± 0.04
	Walker2d	0.23 ± 0.04	0.12 ± 0.00	0.09 ± 0.06	0.11 ± 0.00	0.05 ± 0.04	0.32 ± 0.16	<0.01 ± 0.00
	HalfCheetah	0.36 ± 0.00	0.37 ± 0.00	0.03 ± 0.01	<0.01 ± 0.00	0.32 ± 0.03	0.10 ± 0.00	0.08 ± 0.01
	Pendulum	0.03 ± 0.03	0.08 ± 0.03	0.98 ± 0.00	0.85 ± 0.13	0.07 ± 0.03	0.13 ± 0.00	<0.01 ± 0.01
	Acrobot	0.04 ± 0.01	<0.01 ± 0.00	<0.01 ± 0.00	0.28 ± 0.06	0.10 ± 0.06	0.22 ± 0.06	0.01 ± 0.01

Table 2: Comparison of OPE methods across environments. Error bars represent \pm one standard error across 5 seeds; any regret shown as <0.01 is nonzero but rounds to zero at two decimals.

Baselines We include the following model-free estimators: **Fitted Q-Evaluation (FQE)** [21], **Doubly-Robust OPE (DR)** [36], **Importance Sampling (IS)** [32], and **Density Ratio Estimation (DRE)** [30]. We also include the following model-based estimators: **Model-Based (MB)** [18, 39], and **Policy-Guided Diffusion (PGD)** [15]. The implementation details are provided in Appendix H.

Metrics Each baseline method is evaluated on each pair of behavior dataset and target policy for 5 random seeds. We also generate ground-truth estimates of each target policy value by running each policy in the environment. We evaluate the performance of each baseline using the **Log Root Mean Squared Error (LogRMSE)**, the **Spearman Correlation**, and the **Regret@1** calculated as the difference in return between the best policy selected using the baseline policy value estimates and the actual best policy. Furthermore, to compare metrics consistently across tasks, we normalize the returns following [13]. Appendix I contains the technical details for metric calculation.

4.2 Discussion

Table 2 summarizes the performance of each method per domain, while Figure 2 summarizes the aggregated performance averaged across all domains. STITCH-OPE outperforms all baselines in 11 out of 15 instances (shown in bold), with general agreement among the different metrics. STITCH-OPE soundly outperforms both single-step (MB) and full-trajectory (PGD) model-based methods. This reaffirms our argument in Section 3.3 that intermediate values of w provide a good balance between compositionality and compounding errors. Furthermore, STITCH-OPE performs particularly well according to rank correlation and Regret@1 (with very low standard error) and can accurately rank and identify the best-performing policy. This suggests that the target policy score function (with the negative behavior term) provides very informative guidance during denoising, allowing it to correctly evaluate target policies of varying levels of ability, even as some of those policies deviate significantly from the behavior policy. Finally, we see that STITCH-OPE performance remains consistent across the problem dimension, highlighting the scalability of diffusion when applied to OPE for high-dimensional problems.

		FQE	DRE	MBR	PGD	Ours
Log RMSE ↓	Hopper	-0.21 ± 0.01	-0.38 ± 0.00	-1.56 ± 0.02	-0.89 ± 0.00	-1.65 ± 0.01
	Walker2d	-0.59 ± 0.01	-0.49 ± 0.00	-0.81 ± 0.01	-0.50 ± 0.00	-1.20 ± 0.01
	HalfCheetah	-0.19 ± 0.00	-1.19 ± 0.00	-0.24 ± 0.01	-0.96 ± 0.00	-0.50 ± 0.00
Rank Corr. ↑	Hopper	0.35 ± 0.06	0.35 ± 0.04	0.68 ± 0.02	0.45 ± 0.00	0.81 ± 0.01
	Walker2d	0.03 ± 0.04	0.45 ± 0.03	0.47 ± 0.02	0.52 ± 0.01	0.46 ± 0.09
	HalfCheetah	0.59 ± 0.01	0.80 ± 0.03	0.75 ± 0.05	0.46 ± 0.06	0.81 ± 0.02
Regret@1 ↓	Hopper	0.06 ± 0.03	0.41 ± 0.22	0.18 ± 0.00	<0.01 ± 0.00	<0.01 ± 0.00
	Walker2d	0.24 ± 0.02	0.59 ± 0.13	0.17 ± 0.02	0.23 ± 0.00	0.03 ± 0.00
	HalfCheetah	<0.01 ± 0.00	<0.01 ± 0.00	0.03 ± 0.01	0.02 ± 0.01	0.02 ± 0.01

Table 3: Comparison of OPE methods across environments when the target policy is a diffusion policy; any regret shown as <0.01 is nonzero but rounds to zero at two decimals.

4.3 Off-Policy Evaluation with Diffusion Policies

To demonstrate the ability of STITCH-OPE to evaluate more complex policy classes, we replace target policies with diffusion policies, which have led to significant advances in robotics [7, 40] (see Appendix K for details). Since STITCH-OPE only requires the score of the target policy, it is computationally straightforward to perform OPE with diffusion policies, which is not the case for other estimators that require an explicit probability distribution $\pi_i(a|s)$ over actions (i.e. IS, DR). D4RL results are provided in Table 3. We see that STITCH-OPE outperforms all other baselines in 6 out of 9 instances, demonstrating robust OPE performance across multiple target policy classes.

4.4 Ablations

We conduct additional experiments to test the sensitivity of STITCH-OPE to the choice of guidance coefficients (α and λ) and the window size w . Due to space limitations, we defer results to Appendix L. In summary, the best performance occurs when $0 < \lambda < \alpha$, reaffirming our claim in Section 3.2 that optimal regularization occurs for small λ . The best performance also occurs for $w = 8$, showing that STITCH-OPE provides an optimal balance between autoregressive and full-trajectory diffusion.

5 Limitations

The theory of STITCH-OPE relies on (standard) Assumptions 3.1 and 3.2. In practice, if there exist many (s, a) pairs such that $\beta(a|s) = 0$ but $\pi(a|s) > 0$, then the behavior data may be incomplete and diffusion guidance could generate infeasible trajectories and produce biased estimates of $J(\pi)$. Diffusion models trained on image data in other applications are often easy to interpret; however, evaluating the fidelity of the trajectories generated from the trained behavior model $p_\beta(\tau)$ is challenging when the state is complex and difficult to interpret or partially observable. Finally, while STITCH-OPE has demonstrated excellent performance on existing OPE benchmarks, it remains unanswered whether its benefits also apply to domain-specific problems outside robotics.

6 Conclusion

We presented STITCH-OPE for off-policy evaluation in high-dimensional, long-horizon environments. STITCH-OPE trains a conditional diffusion model to generate behavior sub-trajectories, and applies diffusion guidance using the score of the target policy to correct the distribution shift induced by the target policy. Our novelties include trajectory stitching and negative behavior policy guidance, which were shown to improve composition and generalization. Using D4RL and OpenAI Gym benchmarks, we showed that STITCH-OPE outperforms state-of-the-art OPE methods across MSE, correlation and regret metrics. Future work could investigate online data collection to address severe distribution shift, or explore ways to adapt the guidance coefficients or incorporate additional knowledge into the guidance function (e.g. additional structure on the dynamics). It also remains an open question whether the advantages of STITCH-OPE apply to offline policy optimization.

References

- [1] A. Ajay, Y. Du, A. Gupta, J. B. Tenenbaum, T. S. Jaakkola, and P. Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=sP1fo2K9DFG>.
- [2] P. Alquier and J. Ridgway. Concentration of tempered posteriors and of their variational approximations. *The Annals of Statistics*, 48(3):1475–1497, 2020.
- [3] D. Brandfonbrener, W. Whitney, R. Ranganath, and J. Bruna. Offline rl without off-policy evaluation. *Advances in neural information processing systems*, 34:4933–4946, 2021.
- [4] G. Brockman. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [5] F. Cérou, P. Héas, and M. Rousset. Adaptive reduced tempering for bayesian inverse problems and rare event simulation. *arXiv preprint arXiv:2410.18833*, 2024.
- [6] B. Chen, D. M. Monsó, Y. Du, M. Simchowitz, R. Tedrake, and V. Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [7] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023.
- [8] S. Dankwa and W. Zheng. Twin-delayed ddpg: A deep reinforcement learning technique to model a continuous movement of an intelligent robot agent. In *Proceedings of the 3rd international conference on vision, image and signal processing*, pages 1–5, 2019.
- [9] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [10] M. Farajtabar, Y. Chow, and M. Ghavamzadeh. More robust doubly robust off-policy evaluation. In *International Conference on Machine Learning*, pages 1447–1456. PMLR, 2018.
- [11] M. Fard and J. Pineau. Pac-bayesian model selection for reinforcement learning. *Advances in Neural Information Processing Systems*, 23, 2010.
- [12] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- [13] J. Fu, M. Norouzi, O. Nachum, G. Tucker, ziyu wang, A. Novikov, M. Yang, M. R. Zhang, Y. Chen, A. Kumar, C. Paduraru, S. Levine, and T. Paine. Benchmarks for deep off-policy evaluation. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=kWSeGEeHvF8>.
- [14] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [15] M. T. Jackson, M. Matthews, C. Lu, B. Ellis, S. Whiteson, and J. N. Foerster. Policy-guided diffusion. In *Reinforcement Learning Conference*, 2024.
- [16] M. Janner, Q. Li, and S. Levine. Offline reinforcement learning as one big sequence modeling problem. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 1273–1286. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/099fe6b0b444c23836c4a5d07346082b-Paper.pdf.
- [17] M. Janner, Y. Du, J. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior synthesis. In *International Conference on Machine Learning*, pages 9902–9915. PMLR, 2022.
- [18] N. Jiang and L. Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International conference on machine learning*, pages 652–661. PMLR, 2016.

- [19] N. Jiang, A. Kulesza, S. Singh, and R. Lewis. The dependence of effective planning horizon on model accuracy. In *Proceedings of the 2015 international conference on autonomous agents and multiagent systems*, pages 1181–1189, 2015.
- [20] R. Kidambi, A. Rajeswaran, P. Netrapalli, and T. Joachims. Morel: Model-based offline reinforcement learning. *Advances in neural information processing systems*, 33:21810–21823, 2020.
- [21] H. Le, C. Voloshin, and Y. Yue. Batch policy learning under constraints. In *International Conference on Machine Learning*, pages 3703–3712. PMLR, 2019.
- [22] S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [23] G. Li, L. Shi, Y. Chen, Y. Chi, and Y. Wei. Settling the sample complexity of model-based offline reinforcement learning. *The Annals of Statistics*, 52(1):233–260, 2024.
- [24] Q. Liu, L. Li, Z. Tang, and D. Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. *Advances in neural information processing systems*, 31, 2018.
- [25] V. Liu, P. Nagarajan, A. Patterson, and M. White. When is offline policy selection sample efficient for reinforcement learning? *arXiv preprint arXiv:2312.02355*, 2023.
- [26] Y. Liu, P. L. Bacon, and E. Brunskill. Understanding the curse of horizon in off-policy evaluation via conditional importance sampling. In *International Conference on Machine Learning*, pages 6184–6193. PMLR, 2020.
- [27] Y. Liu, W. Chen, Y. Bai, G. Li, W. Gao, and L. Lin. Aligning cyber space with physical world: A comprehensive survey on embodied ai. *CoRR*, 2024.
- [28] C. Lu, P. J. Ball, Y. W. Teh, and J. Parker-Holder. Synthetic experience replay. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=6jNQ1AY1Uf>.
- [29] L. Mao, H. Xu, X. Zhan, W. Zhang, and A. Zhang. Diffusion-DICE: In-sample diffusion guidance for offline reinforcement learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=EI19qmMmvy>.
- [30] A. Mousavi, L. Li, Q. Liu, and D. Zhou. Black-box off-policy estimation for infinite-horizon reinforcement learning. In *International Conference on Learning Representations*, 2020.
- [31] O. Nachum, Y. Chow, B. Dai, and L. Li. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. *Advances in neural information processing systems*, 32, 2019.
- [32] D. Precup, R. S. Sutton, and S. P. Singh. Eligibility traces for off-policy policy evaluation. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pages 759–766, 2000.
- [33] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [34] M. Rigter, J. Yamada, and I. Posner. World models via policy-guided trajectory diffusion. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=9Ccg00LhKG>.
- [35] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- [36] P. Thomas and E. Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 2139–2148. PMLR, 2016.

- 393 [37] M. Uehara and W. Sun. Pessimistic model-based offline rl: Pac bounds and posterior sampling
394 under partial coverage. *CoRR*, 2021.
- 395 [38] M. Uehara, C. Shi, and N. Kallus. A review of off-policy evaluation in reinforcement learning.
396 *arXiv preprint arXiv:2212.06355*, 2022.
- 397 [39] C. Voloshin, H. M. Le, N. Jiang, and Y. Yue. Empirical study of off-policy policy evaluation for
398 reinforcement learning. In *Thirty-fifth Conference on Neural Information Processing Systems*
399 *Datasets and Benchmarks Track (Round 1)*, 2021.
- 400 [40] Z. Wang, J. J. Hunt, and M. Zhou. Diffusion policies as an expressive policy class for offline
401 reinforcement learning. In *The Eleventh International Conference on Learning Representations*,
402 2023. URL <https://openreview.net/forum?id=AHvFDPi-FA>.
- 403 [41] T. Xie, Y. Ma, and Y.-X. Wang. Towards optimal off-policy evaluation for reinforcement
404 learning with marginalized importance sampling. *Advances in neural information processing*
405 *systems*, 32, 2019.
- 406 [42] M. Yang, O. Nachum, B. Dai, L. Li, and D. Schuurmans. Off-policy evaluation via the
407 regularized lagrangian. *Advances in Neural Information Processing Systems*, 33:6551–6561,
408 2020.
- 409 [43] M. Yang, B. Dai, O. Nachum, G. Tucker, and D. Schuurmans. Offline policy selection under
410 uncertainty. In *International Conference on Artificial Intelligence and Statistics*, pages 4376–
411 4396. PMLR, 2022.
- 412 [44] T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma. Mopo: Model-
413 based offline policy optimization. *Advances in Neural Information Processing Systems*, 33:
414 14129–14142, 2020.
- 415 [45] M. R. Zhang, T. Paine, O. Nachum, C. Paduraru, G. Tucker, ziyu wang, and M. Norouzi.
416 Autoregressive dynamics models for offline policy evaluation and optimization. In *International*
417 *Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=kmqjgSNXby)
418 [id=kmqjgSNXby](https://openreview.net/forum?id=kmqjgSNXby).
- 419 [46] R. Zhang, B. Dai, L. Li, and D. Schuurmans. Gendice: Generalized offline estimation of
420 stationary values. In *International Conference on Learning Representations*, 2020.
- 421 [47] Z. Zhu, H. Zhao, H. He, Y. Zhong, S. Zhang, H. Guo, T. Chen, and W. Zhang. Diffusion models
422 for reinforcement learning: A survey. *arXiv preprint arXiv:2311.01223*, 2023.

1. Claims

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

Answer: [Yes]

Justification: All theoretical and empirical claims have been summarized in the abstract and match the results attained in the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

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

Answer: [Yes]

Justification: Section 5 discusses the assumptions required to obtain the theoretical bound of the method and attain good empirical performance. We also mention that, while the approach performs well on standard benchmark problem sets, its performance on domain-specific problems is unclear.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

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: Appendix C and Appendix D contain all related definitions, theorems and proofs that support the theoretical results stated in the main paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

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: Appendix H discusses how baseline methods were setup as well as all hyper-parameters needed to reproduce the results. Appendix G discusses policy calculation. Appendix F discusses how the domains were set up, and Appendix I discusses how metrics were calculated. Appendix J discusses hyper-parameters needed to run STITCH-OPE on all problems. Together, these sections were included to allow reproduction of the experiments to the best of our ability.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the 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: [Yes]

Justification: Anonymized code is included in the zip file as part of the supplementary material, along with instructions to run the code in a readme file.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

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: The Appendix contains all relevant training and test details, hyper-parameters and other important design considerations.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

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 tables of numerical results and plots clearly show standard error bars and intervals where appropriate.

Guidelines:

- The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

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: Appendix M provides the computer resources used and the total time of experiments of our method.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

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: We have read and adhere to the code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

Answer: [NA]

Justification: We believe the current paper is foundational research, thus we do not foresee any direct societal impacts of the work.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

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: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

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 current work uses code bases from other packages. Their accompanying papers were cited, URL links to the codebases were included, and their licenses (where available) were mentioned in the Appendix.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.

- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

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

Answer: [\[Yes\]](#)

Justification: We provide code to reproduce the experiments as part of the supplementary zip file. Instructions for running the code are provided in a readme file included in the code.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

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: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

746

16. **Declaration of LLM usage**

747

748

749

750

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.

751

Answer: [NA]

752

Justification: The research method does not involve LLMs.

753

Guidelines:

- 754
- 755
- 756
- 757
- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
 - Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

758
759
760
761
762

STITCH-OPE: Trajectory Stitching with Guided Diffusion for Off-Policy Evaluation

Supplementary Material

Abstract

763
764
765
766
767
768

This supplement to the paper discusses algorithmic and experiment details that were not included in the main paper due to space limitations. It includes proofs of all main theoretical claims, as well as all configurations and parameter settings that are required to reproduce the experiments. It includes additional experiments and ablation studies that were excluded from the main paper due to space limitations. It also includes a review of the recent literature on off-policy evaluation in RL.

769

Contents

770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790

A Pedagogical Example for Guided Diffusion	21
B GaussianWorld Domain	21
C Proof that Conditional Diffusion Increases Entropy	22
D Theoretical Analysis	23
D.1 Assumptions and Definitions	23
D.2 Analysis of the Bias	24
D.3 Analysis of the Variance	26
D.4 Proof of the Bias-Variance Decomposition (Theorem 3.3)	30
E Pseudocode	30
F Domains	30
G Policies	31
H Baselines	32
I Metrics	34
J STITCH-OPE Training and Hyper-Parameter Details	35
K Diffusion Policy Training and Evaluation	36
L Additional Experiments	36
L.1 Sensitivity to Guidance Coefficients α and λ	36
L.2 Sensitivity to Window Size w	37
L.3 Trajectory Visualizations	37
M Computing Resources	38
N Related Work	38

791 A Pedagogical Example for Guided Diffusion

792 We consider the simple two-component mixture of Gaussians with density function

$$p(x) = 0.5\mathcal{N}(x; 1, 0.5^2) + 0.5\mathcal{N}(x; -1, 0.5^2),$$

793 where $\mathcal{N}(x; \mu, \sigma^2)$ is the density function of a $\mathcal{N}(\mu, \sigma^2)$ distribution. Using the standard substitution

$$\epsilon(x^k, k) = -\sigma_k \nabla \log p(x^k)$$

794 in the backward diffusion process, produces the backward diffusion process

$$x^{k-1}|x^k \sim \mathcal{N}((x^k + (1 - \alpha_k)\nabla \log p(x^k))/\sqrt{\alpha_k}, \sigma_k^2).$$

795 To illustrate the effects of a guidance function on the sampling process, we consider the guidance
796 function associated with the (unscaled) score of a $\mathcal{N}(1, 0.5^2)$ distribution, i.e.

$$g(x) = -(x - 1)/0.5^2.$$

797 Then, the guided backward diffusion process has mean:

$$\begin{aligned} & \frac{x^k + (1 - \alpha_k)\nabla \log p(x^k)}{\sqrt{\alpha_k}} + \sigma_k^2 g(x^k) \\ &= \frac{x^k + (1 - \alpha_k)(\nabla \log p(x^k) + \sigma_k^2 g(x^k)\sqrt{\alpha_k}/(1 - \alpha_k))}{\sqrt{\alpha_k}}, \end{aligned}$$

798 which corresponds to a standard backward diffusion process with the modified score function

$$\nabla \log p(x^k) + \sigma_k^2 g(x^k)\sqrt{\alpha_k}/(1 - \alpha_k),$$

799 which would place more weight on the rightmost mode of the Gaussian mixture during the backward
800 diffusion process.

801 We run the backward denoising diffusion process using the exact score function $\nabla \log p(x^k)$. The
802 sampling distributions of x^k are plotted at various denoising time steps k in Figure 3.

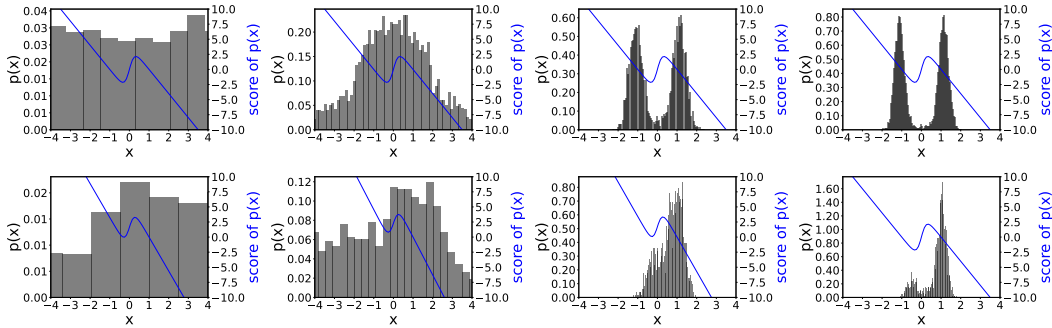


Figure 3: Pedagogical example illustrating guided diffusion sample generation for a Gaussian mixture $0.5\mathcal{N}(1, 0.5^2) + 0.5\mathcal{N}(-1, 0.5^2)$. **Top row:** histograms of samples from unguided backward diffusion at steps $k = 8, 6, 4, 0$, where $\nabla \log p(x)$ is the score of the Gaussian mixture shown in blue. **Bottom row:** histograms of samples from guided diffusion (2) using the score function of a $\mathcal{N}(1, 0.5^2)$ distribution, i.e. $g(x) = -(x - 1)/0.5^2$. The modified score function corresponding to the guided diffusion process is shown in blue. The guided score function (the score of the actual sampling density) is significantly shifted and skewed, relative to the original score function, at the intermediate denoising time steps ($k = 6, 4$). This ensures that the right mode of the Gaussian mixture is sampled more frequently during denoising.

803 B GaussianWorld Domain

804 The GaussianWorld domain is a toy 2-dimensional Markov decision process defined designed to
805 illustrate and compare generalization and compositionality of diffusion models (Table 1). It is defined
806 as follows:

807 **Decision Epochs** The decision epochs are $t = 0, 1, 2, \dots T$ where we set $T = 128$ in our experi-
 808 ments.

809 **State Space** $\mathcal{S} = \mathbb{R}^2$ describes all positions (x_t, y_t) of a particle in space at every decision epoch t .
 810 It is assumed that x_t is the x-coordinate and y_t is the y-coordinate. The initial state is $s_0 = (0, 0)$
 811 unless otherwise specified.

812 **Action Space** $\mathcal{A} = \mathbb{R}$ describes the (counterclockwise) angle of the movement vector of the particle
 813 at every decision epoch, relative to the horizontal.

814 **Transitions** Letting a_t be the angle of movement of the particle at time t , the transitions of x_t and
 815 y_t are defined as follows:

$$x_{t+1} = x_t + 0.02 \cdot \cos(a_t + \varepsilon_t), \quad y_{t+1} = y_t + 0.02 \cdot \sin(a_t + \varepsilon_t), \quad \varepsilon_t \sim \mathcal{N}(0, 0.2^2).$$

816 Here, ε_t is an i.i.d. Gaussian noise added to the actions before they are applied by the controller.

817 **Reward Function and Discount** The problem is not solved so we leave the reward unspecified.
 818 We also leave the discount factor unspecified.

819 C Proof that Conditional Diffusion Increases Entropy

820 We begin with the following definitions.

821 **Definition C.1** (Entropy). Let $p(x)$ be a density function of a random variable X with support \mathcal{X} .
 822 The *entropy* of X is defined as

$$H(X) = \int_{\mathcal{X}} p(x) \log \left(\frac{1}{p(x)} \right) dx.$$

823 **Definition C.2** (Conditional Entropy). The *conditional entropy* of X given Y on support \mathcal{Y} is defined
 824 as

$$H(X|Y) = \mathbb{E}_{y \in \mathcal{Y}} [H(X|Y = y)].$$

825 Our goal is to prove

826 **Theorem C.3.** Let S_t be the random state at time t sampled according to the conditional distribution
 827 $p(S_{t+1} = s | S_t = x, A_t = u)$, and let A_t be a random action following some conditional distribution
 828 $p(A_t = a | S_t = x)$. Then $H(\tau | S_t) \geq H(\tau | S_0, A_0 \dots S_t)$, where τ is a (random) sub-trajectory
 829 beginning in state S_t .

830 *Proof.* First, letting $U = (S_0, A_0, \dots S_{t-1}, A_{t-1})$, observe that:

$$\begin{aligned} & H(U, \tau | S_t = s) \\ &= \iint p(U = u, \tau | S_t = s) \log \left(\frac{1}{p(U = u, \tau | S_t = s)} \right) du d\tau \\ &= \iint p(U = u, \tau | S_t = s) \log \left(\frac{1}{p(U = u | S_t = s) p(\tau | U = u, S_t = s)} \right) du d\tau \\ &= \int p(U = u | S_t = s) \log \left(\frac{1}{p(U = u | S_t = s)} \right) du \\ &\quad + \int p(U = u | S_t = s) \int p(\tau | U = u, S_t = s) \log \left(\frac{1}{p(\tau | U = u, S_t = s)} \right) du d\tau \\ &= H(U | S_t = s) + \mathbb{E}_{u \in \mathcal{U} | S_t = s} [H(\tau | U = u, S_t = s)]. \end{aligned}$$

831 Next, using the additivity property of expectation and law of total expectation:

$$\begin{aligned} H(U, \tau | S_t) &= \mathbb{E}_{s \in \mathcal{S}_t} [H(U | S_t = s)] + \mathbb{E}_{s \in \mathcal{S}_t, u \in \mathcal{U}} [H(\tau | U = u, S_t = s)] \\ &= H(U | S_t) + H(\tau | U, S_t). \end{aligned}$$

832 Next, we prove sub-additivity of conditional entropy:

$$\begin{aligned}
& H(U, \tau | S_t = s) - H(U | S_t = s) - H(\tau | S_t = s) \\
&= \iint p(U = u, \tau | S_t = s) \log \left(\frac{1}{p(U = u, \tau | S_t = s)} \right) du d\tau \\
&\quad - \int p(U = u | S_t = s) \log \left(\frac{1}{p(U = u | S_t = s)} \right) du - \int p(\tau | S_t = s) \log \left(\frac{1}{p(\tau | S_t = s)} \right) d\tau \\
&= \iint p(U = u, \tau | S_t = s) \log \left(\frac{p(\tau | S_t = s) p(U = u | S_t = s)}{p(U = u, \tau | S_t = s)} \right) du d\tau \\
&\leq \log \iint p(U = u, \tau | S_t = s) \left(\frac{p(\tau | S_t = s) p(U = u | S_t = s)}{p(U = u, \tau | S_t = s)} \right) du d\tau \\
&= \log 1 = 0,
\end{aligned}$$

833 where the inequality in the derivation follows by Jensen's inequality. This implies that

$$H(U, \tau | S_t = s) \leq H(U | S_t = s) + H(\tau | S_t = s).$$

834 Taking expectation of both sides with respect to S_t , and using the monotonicity and additivity
835 properties of expectation:

$$\begin{aligned}
H(U, \tau | S_t) &= \mathbb{E}_{s \in S_t} [H(U, \tau | S_t = s)] \\
&\leq \mathbb{E}_{s \in S_t} [H(U | S_t = s) + H(\tau | S_t = s)] = H(U | S_t) + H(\tau | S_t).
\end{aligned}$$

836 Finally, putting it all together:

$$H(\tau | U, S_t) = H(U, \tau | S_t) - H(U | S_t) \leq H(U | S_t) + H(\tau | S_t) - H(U | S_t) = H(\tau | S_t),$$

837 which completes the proof. \square

838 D Theoretical Analysis

839 D.1 Assumptions and Definitions

840 We decompose a full trajectory of length T into $N = T/w$ non-overlapping sub-trajectories (or
841 chunks), each of length w . Each *chunk* $S_i \in \mathcal{T}^{(w)}$ is defined as

$$S_i := (s_{iw}, a_{iw}, s_{iw+1}, a_{iw+1}, \dots, s_{(i+1)w}).$$

842 Let the *full trajectory* be defined as

$$S = (S_0, S_1, \dots, S_{N-1}).$$

843 We define the *boundary state* X_i as the initial state of chunk S_i :

$$X_i := s_{iw}, \quad i = 0, 1 \dots N,$$

844 which form the backbone of the generative process.

845 We assume the following factored generative process for trajectories

$$p(S_0, S_1, \dots, S_{N-1}) = p(X_0) \prod_{i=0}^{N-1} p(S_i | X_i) p(X_{i+1} | S_i).$$

846 This implies that the boundary state sequence $X = (X_0, X_1, \dots, X_N)$ forms a first-order Markov
847 chain

$$p(X_{i+1} | S_i) = p(X_{i+1} | X_i).$$

848 Each chunk S_i produces a scalar discounted return Y_i , defined as

$$Y_i := f(S_i) = \sum_{j=0}^{w-1} \gamma^j \hat{R}(s_{iw+j}, a_{iw+j}),$$

849 where \hat{R} is a learned reward model, and $\gamma \in [0, 1]$ is the discount factor.

850 Given a bound $R_{\max} < \infty$ on the absolute reward, we define the *maximum per-chunk return bound*
 851 as:

$$B_w := \sum_{j=0}^{w-1} \gamma^j R_{\max} = \frac{R_{\max}(1 - \gamma^w)}{1 - \gamma} \Rightarrow |Y_i| \leq B_w.$$

852 The cumulative return over the full trajectory is approximated by

$$\hat{J} = \sum_{i=0}^{N-1} \gamma^{iw} Y_i,$$

853 and the expected return under the target policy π is:

$$J(\pi) := \mathbb{E}_{p_\pi}[\hat{J}] = \mathbb{E}_{p_\pi} \left[\sum_{i=0}^{N-1} \gamma^{iw} Y_i \right].$$

854 **Definition D.1** (Chunked Behavior Distributions). Let $p_\beta^{(w)}$ denote the true distribution over behavior
 855 chunks S_i , and let $\hat{p}_\beta^{(w)}$ be the learned conditional distribution modeled by the diffusion process.
 856 These distributions describe how chunks are generated given boundary states:

$$p_\beta^{(w)}(S_i | X_i), \quad \hat{p}_\beta^{(w)}(S_i | X_i).$$

857 **Definition D.2** (Total Variation Distance). The *total variation distance* between two probability
 858 distributions P and Q over the same measurable space \mathcal{X} is defined as

$$\text{TV}(P, Q) := \sup_{A \subseteq \mathcal{X}} |P(A) - Q(A)|.$$

859 We now restate the two assumptions presented in the main text for convenience.

860 **Assumption D.3** (Bounded Likelihood Ratio). There is a constant κ such that $\frac{\pi(a|s)}{\beta(a|s)} \leq \kappa$ for all
 861 $s \in \mathcal{S}$ and $a \in \mathcal{A}$.

862 Note that this assumption can be easily verified in our experimental setting. Since the action spaces
 863 are closed intervals and the behavior and target policy distributions are both represented as truncated
 864 Gaussian distributions, the ratio of the two policies is bounded over the action space.

865 **Assumption D.4** (Chunk-wise Model Fit). The total variation distance between the true chunk
 866 distribution $p_\beta^{(w)}$ and the learned conditional distribution $\hat{p}_\beta^{(w)}$ is bounded by some constant $\delta_\beta > 0$,

$$\text{TV}(p_\beta^{(w)}, \hat{p}_\beta^{(w)}) \leq \delta_\beta.$$

867 D.2 Analysis of the Bias

868 We begin by bounding the total variation distance between the true target distribution $p_\pi^{(w)}$ and the
 869 guided model $\hat{p}_\pi^{(w)}$.

870 **Lemma D.5.** *The total variation distance between the guided model $\hat{p}_\pi^{(w)}$ and the true target*
 871 *distribution $p_\pi^{(w)}$ satisfies*

$$\text{TV}(p_\pi^{(w)}, \hat{p}_\pi^{(w)}) \leq \kappa^2 \cdot \delta_\beta$$

872 *Proof.* By the definition of total variation distance

$$\text{TV}(p_\pi^{(w)}, \hat{p}_\pi^{(w)}) = \frac{1}{2} \int \left| p_\pi^{(w)}(\tau) - \hat{p}_\pi^{(w)}(\tau) \right| d\tau.$$

873 Using the reweighted form of each distribution

$$\text{TV}(p_\pi^{(w)}, \hat{p}_\pi^{(w)}) = \frac{1}{2} \int \left| \left(p_\beta^{(w)}(\tau) - \hat{p}_\beta^{(w)}(\tau) \right) \cdot \prod_{j=0}^{w-1} \frac{\pi(a_j | s_j)}{\beta(a_j | s_j)} \right| d\tau,$$

874 and applying the bound on the likelihood ratio (Assumption D.3):

$$\text{TV}(p_\pi^{(w)}, \hat{p}_\pi^{(w)}) \leq \frac{\kappa^w}{2} \int \left| p_\beta^{(w)}(\tau) - \hat{p}_\beta^{(w)}(\tau) \right| d\tau = \kappa^w \cdot \text{TV}(p_\beta^{(w)}, \hat{p}_\beta^{(w)}) \leq \kappa^w \cdot \delta_\beta.$$

875 This completes the proof. \square

876 Let the total variation distance between the true target distribution and the guided diffusion model be
877 denoted by

$$\delta_\pi := \text{TV}\left(p_\pi^{(w)}, \hat{p}_\pi^{(w)}\right).$$

878 By Lemma D.5, we have the bound

$$\delta_\pi \leq \kappa^w \cdot \delta_\beta.$$

879 We now derive a bound on the absolute bias of the estimated return when sampling chunks from the
880 guided model $\hat{p}_\pi^{(w)}$ instead of the true target distribution $p_\pi^{(w)}$.

881 **Lemma D.6** (Expectation Difference Bound via Total Variation). *Let p and q be two probability
882 densities on a probability space \mathcal{X} . Let*

$$\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)|$$

883 *be the supremum norm of a bounded function $f : \mathcal{X} \rightarrow \mathbb{R}$, and let:*

$$\|p - q\|_1 = \int_{\mathcal{X}} |p(x) - q(x)| dx, \quad \text{TV}(p, q) = \frac{1}{2} \|p - q\|_1.$$

884 *Then*

$$|\mathbb{E}_{x \sim p}[f(x)] - \mathbb{E}_{x \sim q}[f(x)]| \leq 2 \|f\|_\infty \text{TV}(p, q).$$

Proof.

$$\begin{aligned} |\mathbb{E}_p[f] - \mathbb{E}_q[f]| &= \left| \int_{\mathcal{X}} f(x) p(x) dx - \int_{\mathcal{X}} f(x) q(x) dx \right| \\ &= \left| \int_{\mathcal{X}} f(x) (p(x) - q(x)) dx \right| \leq \int_{\mathcal{X}} |f(x)| |p(x) - q(x)| dx \\ &\leq \|f\|_\infty \int_{\mathcal{X}} |p(x) - q(x)| dx = 2 \|f\|_\infty \text{TV}(p, q). \end{aligned}$$

885 This completes the proof. \square

886 **Lemma D.7** (Marginal TV Bound via Conditional TV). *Let $p(x | s)$ and $\hat{p}(x | s)$ be conditional
887 densities over chunk $x \in \mathcal{T}^{(w)}$, given state $s \in \mathcal{S}$, and let $\mu(s)$ denote the marginal distribution over
888 s . Then*

$$\text{TV}\left(\int p(x | s) \mu(s) ds, \int \hat{p}(x | s) \mu(s) ds\right) \leq \int \text{TV}(p(\cdot | s), \hat{p}(\cdot | s)) \mu(s) ds.$$

889 *In particular, if $\text{TV}(p(\cdot | s), \hat{p}(\cdot | s)) \leq \epsilon$ for all s , then*

$$\text{TV}(p, \hat{p}) \leq \epsilon.$$

890 *Proof.* Let $p(x) = \int p(x | s) \mu(s) ds$, $\hat{p}(x) = \int \hat{p}(x | s) \mu(s) ds$. Then:

$$\begin{aligned} \text{TV}(p, \hat{p}) &= \frac{1}{2} \int |p(x) - \hat{p}(x)| dx \\ &= \frac{1}{2} \int \left| \int \mu(s) [p(x | s) - \hat{p}(x | s)] ds \right| dx \\ &\leq \frac{1}{2} \iint \mu(s) |p(x | s) - \hat{p}(x | s)| ds dx \quad (\text{by Jensen's inequality}) \\ &= \int \mu(s) \left[\frac{1}{2} \int |p(x | s) - \hat{p}(x | s)| dx \right] ds \\ &= \int \mu(s) \cdot \text{TV}(p(\cdot | s), \hat{p}(\cdot | s)) ds. \end{aligned}$$

891 If $\text{TV}(p(\cdot | s), \hat{p}(\cdot | s)) \leq \epsilon$ uniformly, the integral is bounded by ϵ . \square

892 **Theorem D.8** (Bias Bound for STITCH-OPE). *The bias of the return estimate under the guided*
 893 *diffusion model satisfies*

$$\left| \mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi) \right| \leq \frac{2B_w}{1 - \gamma^w} \cdot \delta_\pi.$$

894 *Proof.* The return estimator is:

$$\hat{J} = \sum_{i=0}^{N-1} \gamma^{iw} Y_i, \quad \text{where } Y_i = f(S_i) = \sum_{j=0}^{w-1} \gamma^j \hat{R}(s_{iw+j}, a_{iw+j}).$$

895 Thus, the bias is:

$$\left| \mathbb{E}_{\hat{p}_\pi}[\hat{J}] - \mathbb{E}_{p_\pi}[\hat{J}] \right| = \left| \sum_{i=0}^{N-1} \gamma^{iw} (\mathbb{E}_{\hat{p}_\pi}[Y_i] - \mathbb{E}_{p_\pi}[Y_i]) \right| \leq \sum_{i=0}^{N-1} \gamma^{iw} |\mathbb{E}_{\hat{p}_\pi}[Y_i] - \mathbb{E}_{p_\pi}[Y_i]|.$$

896 For each chunk i , Y_i depends only on S_i , with marginal distributions $\hat{p}_\pi^{(w,i)}$ and $p_\pi^{(w,i)}$ under \hat{p}_π and
 897 p_π , respectively. By Lemma D.6 and Lemma D.7

$$|\mathbb{E}_{\hat{p}_\pi}[Y_i] - \mathbb{E}_{p_\pi}[Y_i]| \leq 2 \cdot \sup |Y_i| \cdot \text{TV}(p_\pi^{(w,i)}, \hat{p}_\pi^{(w,i)}).$$

898 Since $|\hat{R}(s, a)| \leq R_{\max}$, the per-chunk return is bounded:

$$|Y_i| \leq \sum_{j=0}^{w-1} \gamma^j R_{\max} = R_{\max} \cdot \frac{1 - \gamma^w}{1 - \gamma}.$$

899 Using Lemma D.5, we know that $\text{TV}(p_\pi^{(w,i)}, \hat{p}_\pi^{(w,i)}) \leq \delta_\pi$. Thus we have

$$|\mathbb{E}_{\hat{p}_\pi}[Y_i] - \mathbb{E}_{p_\pi}[Y_i]| \leq 2 \cdot \frac{R_{\max}(1 - \gamma^w)}{1 - \gamma} \cdot \delta_\pi.$$

900 Summing over chunks:

$$\left| \mathbb{E}_{\hat{p}_\pi}[\hat{J}] - \mathbb{E}_{p_\pi}[\hat{J}] \right| \leq \sum_{i=0}^{N-1} \gamma^{iw} \cdot 2 \cdot \frac{R_{\max}(1 - \gamma^w)}{1 - \gamma} \cdot \delta_\pi = 2 \cdot \frac{R_{\max}(1 - \gamma^w)}{1 - \gamma} \cdot \delta_\pi \cdot \sum_{i=0}^{N-1} \gamma^{iw}.$$

901 The geometric sum is:

$$\sum_{i=0}^{N-1} \gamma^{iw} \leq \sum_{i=0}^{\infty} \gamma^{iw} = \frac{1}{1 - \gamma^w}.$$

902 Thus:

$$\left| \mathbb{E}_{\hat{p}_\pi}[\hat{J}] - \mathbb{E}_{p_\pi}[\hat{J}] \right| \leq 2 \cdot \frac{R_{\max}(1 - \gamma^w)}{1 - \gamma} \cdot \delta_\pi \cdot \frac{1}{1 - \gamma^w} = \frac{2B_w}{1 - \gamma^w} \cdot \delta_\pi.$$

903 This completes the proof. \square

904 **Corollary D.9** (Bias Bound in Terms of Model Fit δ_β). *Under the assumptions*
 905 $\sup_i \text{TV}(p_\pi^{(w,i)}, \hat{p}_\pi^{(w,i)}) \leq \delta_\pi \leq \kappa^w \cdot \delta_\beta$ and $\sup_\tau |\hat{J}(\tau)| \leq \frac{R_{\max}}{1 - \gamma}$, the bias satisfies

$$\left| \mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi) \right| \leq \frac{2B_w}{1 - \gamma^w} \cdot \kappa^w \cdot \delta_\beta.$$

906 D.3 Analysis of the Variance

907 **Lemma D.10** (Conditional Independence of Chunk Rewards). *Let $X_i := s_{iw}$ be the boundary state*
 908 *at the start of chunk S_i , and define:*

$$Y_i := f(S_i) = \sum_{j=0}^{w-1} \gamma^j \hat{R}(s_{iw+j}, a_{iw+j}).$$

909 Assume the generative process satisfies the following properties:

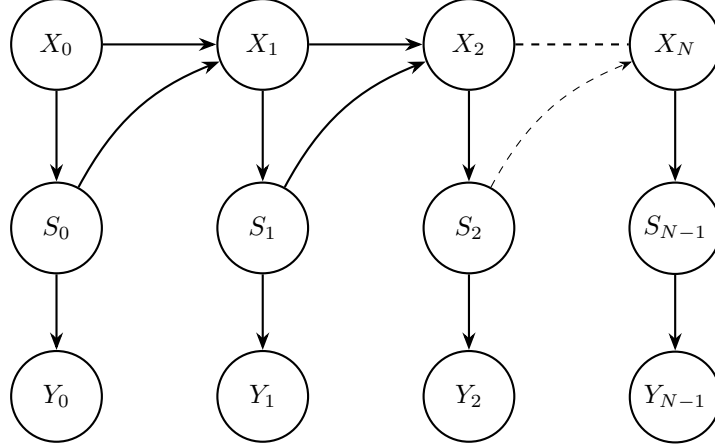


Figure 4: Illustration of the sub-trajectory decomposition. Each chunk S_i generates a reward sequence Y_i and leads to a boundary state X_{i+1} .

910 • Each chunk S_i is generated independently given X_i

911 • The return Y_i is a deterministic function of S_i .

912 Then for all $i \neq j$, the returns Y_i and Y_j are conditionally independent given the full boundary state
913 chain X_0, X_1, \dots, X_N ,

$$Y_i \perp\!\!\!\perp Y_j \mid X_0, \dots, X_N.$$

914 *Proof.* Refer to the graphical model in Figure 4. The nodes X_0, X_1, \dots, X_N form a Markov chain.
915 Each chunk S_i is a child of X_i , and each return Y_i is a child of S_i .

916 Now consider any path from Y_i to Y_j . Such a path must go through:

$$Y_i \leftarrow S_i \leftarrow X_i \rightsquigarrow X_{i+1} \rightsquigarrow \dots \rightsquigarrow X_j \rightarrow S_j \rightarrow Y_j.$$

917 All such paths must traverse through at least one boundary node X_k . Since we are conditioning on all
918 X_0, \dots, X_N , and these nodes are non-colliders on every path from Y_i to Y_j , all such paths are blocked.
919 By the criterion of d-separation (see, e.g. Chapter 8 in [48]), this implies $Y_i \perp\!\!\!\perp Y_j \mid X_0, \dots, X_N$. \square

920 **Theorem D.11** (Variance Bound). *Let \hat{p}_π denote the trajectory distribution induced by the guided
921 diffusion model, and p_π the true trajectory distribution under the target policy. Let \hat{J} be the return
922 estimator using a learned reward model. Then*

$$\text{Var}_{\hat{p}_\pi}(\hat{J}) \leq \text{Var}_{p_\pi}(J) + 10 \left(\frac{T}{w} \right)^2 B_w^2 \kappa^w \delta_\beta + \frac{2B_w^2}{1 - \gamma^{2w}} \kappa^w \delta_\beta,$$

923 where B_w denotes the maximum per-chunk discounted return.

924 *Proof.* We begin by applying the law of total variance under the guided model distribution \hat{p}_π

$$\text{Var}_{\hat{p}_\pi}(\hat{J}) = \mathbb{E}_{\hat{p}_\pi} \left[\text{Var}_{\hat{p}_\pi}(\hat{J} \mid X) \right] + \text{Var}_{\hat{p}_\pi} \left(\mathbb{E}_{\hat{p}_\pi}[\hat{J} \mid X] \right).$$

925 Using Lemma D.10 we have that the chunk-level rewards Y_i and Y_j are conditionally independent
926 given the boundary states X_0, X_1, \dots, X_N :

$$Y_i \perp\!\!\!\perp Y_j \mid X_0, X_1, \dots, X_N \quad \text{for all } i \neq j.$$

927 Using this conditional independence, the variance of the total return under \hat{p}_π factorizes:

$$\text{Var}_{\hat{p}_\pi}[\hat{J} \mid X] = \text{Var}_{\hat{p}_\pi} \left[\sum_{i=0}^{N-1} \gamma^{iw} Y_i \mid X \right] = \sum_{i=0}^{N-1} \gamma^{2iw} \cdot \text{Var}_{\hat{p}_\pi}(Y_i \mid X_i).$$

928 To bound the difference in conditional variances, we apply the law of variance

$$\text{Var}(Y_i | X_i) = \mathbb{E}[Y_i^2 | X_i] - (\mathbb{E}[Y_i | X_i])^2.$$

929 Let us define a bound on the per-chunk return magnitude:

$$B_w := \frac{R_{\max}(1 - \gamma^w)}{1 - \gamma} \Rightarrow |Y_i| \leq B_w, \quad Y_i^2 \leq B_w^2.$$

930 Using Lemma D.6 (Expectation Difference Bound via Total Variation), we have

$$|\mathbb{E}_{p_\pi}[f] - \mathbb{E}_{\hat{p}_\pi}[f]| \leq 2\delta_\pi \cdot \|f\|_\infty.$$

931 Applying this with $f = Y_i$ and $f = Y_i^2$, and using the bound $|Y_i| \leq B_w$, we obtain:

$$|\mathbb{E}_{p_\pi}[Y_i] - \mathbb{E}_{\hat{p}_\pi}[Y_i]| \leq 2\delta_\pi B_w, \quad |\mathbb{E}_{p_\pi}[Y_i^2] - \mathbb{E}_{\hat{p}_\pi}[Y_i^2]| \leq 2\delta_\pi B_w^2.$$

932 We analyze the difference in conditional variances:

$$\begin{aligned} & |\text{Var}_{\hat{p}_\pi}(Y_i | X_i) - \text{Var}_{p_\pi}(Y_i | X_i)| \\ &= |\mathbb{E}_{\hat{p}_\pi}[Y_i^2] - \mathbb{E}_{p_\pi}[Y_i^2] - (\mathbb{E}_{\hat{p}_\pi}[Y_i]^2 - \mathbb{E}_{p_\pi}[Y_i]^2)| \\ &\leq |\mathbb{E}_{\hat{p}_\pi}[Y_i^2] - \mathbb{E}_{p_\pi}[Y_i^2]| + |\mathbb{E}_{\hat{p}_\pi}[Y_i]^2 - \mathbb{E}_{p_\pi}[Y_i]^2| \\ &= |\mathbb{E}_{\hat{p}_\pi}[Y_i^2] - \mathbb{E}_{p_\pi}[Y_i^2]| + |\mathbb{E}_{\hat{p}_\pi}[Y_i] - \mathbb{E}_{p_\pi}[Y_i]| \cdot |\mathbb{E}_{\hat{p}_\pi}[Y_i] + \mathbb{E}_{p_\pi}[Y_i]| \\ &\leq 2\delta_\pi B_w^2 + (2\delta_\pi B_w)(2B_w) \\ &= 6\delta_\pi B_w^2. \end{aligned}$$

933 This uses the triangle inequality and the identity $|a^2 - b^2| = |a - b||a + b|$, along with the bounds

934 $|Y_i| \leq B_w$, $\|Y_i\|_\infty^2 \leq B_w^2$, and total variation guarantees from Lemma D.6. Then

$$|\text{Var}_{\hat{p}_\pi}(Y_i | X_i) - \text{Var}_{p_\pi}(Y_i | X_i)| \leq 6\delta_\pi B_w^2.$$

935 We now return to bounding the first term in the law of total variance

$$\mathbb{E}_{\hat{p}_\pi} \left[\text{Var}_{\hat{p}_\pi}(\hat{J} | X) \right] = \mathbb{E}_{\hat{p}_\pi} \left[\sum_{i=0}^{N-1} \gamma^{2iw} \cdot \text{Var}_{\hat{p}_\pi}(Y_i | X_i) \right].$$

936 Using the bound from the previous step

$$\text{Var}_{\hat{p}_\pi}(Y_i | X_i) \leq \text{Var}_{p_\pi}(Y_i | X_i) + 6\delta_\pi B_w^2.$$

937 Taking expectation over \hat{p}_π on both sides

$$\mathbb{E}_{\hat{p}_\pi} [\text{Var}_{\hat{p}_\pi}(Y_i | X_i)] \leq \mathbb{E}_{\hat{p}_\pi} [\text{Var}_{p_\pi}(Y_i | X_i)] + 6\delta_\pi B_w^2.$$

938 Now, using the expectation difference bound from Lemma D.6 again:

$$|\mathbb{E}_{\hat{p}_\pi}[f] - \mathbb{E}_{p_\pi}[f]| \leq 2\delta_\pi \|f\|_\infty, \quad \text{where} \quad f(X_i) := \text{Var}_{p_\pi}(Y_i | X_i) \leq B_w^2.$$

939 So

$$\mathbb{E}_{\hat{p}_\pi} [\text{Var}_{p_\pi}(Y_i | X_i)] \leq \mathbb{E}_{p_\pi} [\text{Var}_{p_\pi}(Y_i | X_i)] + 2\delta_\pi B_w^2.$$

940 Combining both components

$$\mathbb{E}_{\hat{p}_\pi} [\text{Var}_{\hat{p}_\pi}(Y_i | X_i)] \leq \mathbb{E}_{p_\pi} [\text{Var}_{p_\pi}(Y_i | X_i)] + 8\delta_\pi B_w^2.$$

941 Summing across all chunks:

$$\begin{aligned} \mathbb{E}_{\hat{p}_\pi} \left[\text{Var}_{\hat{p}_\pi}(\hat{J} | X) \right] &= \sum_{i=0}^{N-1} \gamma^{2iw} \cdot \mathbb{E}_{\hat{p}_\pi} [\text{Var}_{\hat{p}_\pi}(Y_i | X_i)] \\ &\leq \sum_{i=0}^{N-1} \gamma^{2iw} (\mathbb{E}_{p_\pi} [\text{Var}_{p_\pi}(Y_i | X_i)] + 8\delta_\pi B_w^2). \end{aligned}$$

942 We can split the sum and factor out constants:

$$\mathbb{E}_{\hat{p}_\pi} \left[\text{Var}_{\hat{p}_\pi}(\hat{J} \mid X) \right] = \sum_{i=0}^{N-1} \gamma^{2iw} \cdot \mathbb{E}_{p_\pi} [\text{Var}_{p_\pi}(Y_i \mid X_i)] + 8\delta_\pi B_w^2 \sum_{i=0}^{N-1} \gamma^{2iw}.$$

943 Let us define the chunk-level return variance

$$\mathbb{E}_{p_\pi} \left[\text{Var}_{p_\pi}(\hat{J} \mid X) \right] := \sum_{i=0}^{N-1} \gamma^{2iw} \cdot \mathbb{E}_{p_\pi} [\text{Var}_{p_\pi}(Y_i \mid X_i)].$$

944 Therefore

$$\boxed{\mathbb{E}_{\hat{p}_\pi} \left[\text{Var}_{\hat{p}_\pi}(\hat{J} \mid X) \right] \leq \mathbb{E}_{p_\pi} \left[\text{Var}_{p_\pi}(\hat{J} \mid X) \right] + \frac{8\delta_\pi B_w^2}{1 - \gamma^{2w}}.}$$

945 To complete the law of total variance, we now analyze the second term:

$$\text{Var}_{\hat{p}_\pi} \left(\mathbb{E}_{\hat{p}_\pi}[\hat{J} \mid X] \right) = \text{Var}_{\hat{p}_\pi}(Z_{\hat{p}}), \quad \text{where } Z_{\hat{p}} := \sum_{k=0}^{N-1} g_k(X_k), \quad g_k(x) := \mathbb{E}_{\hat{p}_\pi}[Y_k \mid X_k = x].$$

946 We define the corresponding ideal (true model) version:

$$Z_p := \sum_{k=0}^{N-1} \tilde{g}_k(X_k), \quad \tilde{g}_k(x) := \mathbb{E}_{p_\pi}[Y_k \mid X_k = x].$$

947 Our goal is to bound the variance difference:

$$\Delta_{\text{mean}} := \text{Var}_{\hat{p}_\pi}(Z_{\hat{p}}) - \text{Var}_{p_\pi}(Z_p) = (M_{\hat{p}} - M_p) - (m_{\hat{p}} - m_p)(m_{\hat{p}} + m_p),$$

948 where $M_{\hat{p}} := \mathbb{E}_{\hat{p}_\pi}[Z_{\hat{p}}^2]$, $m_{\hat{p}} := \mathbb{E}_{\hat{p}_\pi}[Z_{\hat{p}}]$, and similarly for M_p, m_p .

949 Insert and subtract a common term:

$$m_{\hat{p}} - m_p = \sum_{k=0}^{N-1} (\mathbb{E}_{\hat{p}_\pi}[g_k(X_k)] - \mathbb{E}_{\hat{p}_\pi}[\tilde{g}_k(X_k)]) + \sum_{k=0}^{N-1} (\mathbb{E}_{\hat{p}_\pi}[\tilde{g}_k(X_k)] - \mathbb{E}_{p_\pi}[\tilde{g}_k(X_k)]).$$

950 Each term is bounded by $2\delta_\pi B_w$, so $|m_{\hat{p}} - m_p| \leq 4N\delta_\pi B_w$.

951 Expand both squares:

$$\begin{aligned} Z_{\hat{p}}^2 &= \sum_{k=0}^{N-1} g_k^2(X_k) + 2 \sum_{0 \leq k < \ell \leq N-1} g_k(X_k)g_\ell(X_\ell), \\ Z_p^2 &= \sum_{k=0}^{N-1} \tilde{g}_k^2(X_k) + 2 \sum_{0 \leq k < \ell \leq N-1} \tilde{g}_k(X_k)\tilde{g}_\ell(X_\ell). \end{aligned}$$

952 Each term (both diagonal and cross terms) is bounded in total variation with sup-norm B_w^2 , yielding

$$|M_{\hat{p}} - M_p| \leq 2N^2\delta_\pi B_w^2.$$

953 From the bound on the means:

$$|m_{\hat{p}}|, |m_p| \leq NB_w \quad \Rightarrow \quad |m_{\hat{p}} + m_p| \leq 2NB_w.$$

954 So, the product term:

$$|(m_{\hat{p}} - m_p)(m_{\hat{p}} + m_p)| \leq (4N\delta_\pi B_w)(2NB_w) = 8N^2\delta_\pi B_w^2.$$

955 Combining both:

$$|\Delta_{\text{mean}}| = |\text{Var}_{\hat{p}_\pi}(Z_{\hat{p}}) - \text{Var}_{p_\pi}(Z_p)| \leq 2N^2\delta_\pi B_w^2 + 8N^2\delta_\pi B_w^2 = 10N^2\delta_\pi B_w^2,$$

956 which yields

$$\boxed{\left| \text{Var}_{\hat{p}_\pi} \left(\mathbb{E}_{\hat{p}_\pi}[\hat{J} \mid X] \right) - \text{Var}_{p_\pi} \left(\mathbb{E}_{p_\pi}[J \mid X] \right) \right| \leq 10 \cdot \frac{T^2}{w^2} \cdot \delta_\pi B_w^2.}$$

957 Combining the two components from the law of total variance, we conclude:

$$\begin{aligned}
\text{Var}_{\hat{p}_\pi}(\hat{J}) &= \mathbb{E}_{\hat{p}_\pi} \left[\text{Var}_{\hat{p}_\pi}(\hat{J} \mid X) \right] + \text{Var}_{\hat{p}_\pi} \left(\mathbb{E}_{\hat{p}_\pi}[\hat{J} \mid X] \right) \\
&\leq \mathbb{E}_{p_\pi} \left[\text{Var}_{p_\pi}(\hat{J} \mid X) \right] + \frac{8\delta_\pi B_w^2}{1 - \gamma^{2w}} + \text{Var}_{p_\pi}(\mathbb{E}_{p_\pi}[J \mid X]) + 10 \left(\frac{T}{w} \right)^2 \delta_\pi B_w^2 \\
&= \text{Var}_{p_\pi}(J) + 10 \left(\frac{T}{w} \right)^2 \delta_\pi B_w^2 + \frac{8\delta_\pi B_w^2}{1 - \gamma^{2w}}.
\end{aligned}$$

958 By Lemma D.5,

$$\boxed{\text{Var}_{\hat{p}_\pi}(\hat{J}) \leq \text{Var}_{p_\pi}(J) + 10 \left(\frac{T}{w} \right)^2 B_w^2 \kappa^w \delta_\beta + \frac{8B_w^2}{1 - \gamma^{2w}} \kappa^w \delta_\beta,}$$

959 and the proof is complete. \square

960 D.4 Proof of the Bias-Variance Decomposition (Theorem 3.3)

961 Finally, we can bound the mean squared error of STITCH-OPE.

962 **Theorem D.12.** *Under Assumption D.3 and D.4, and using the notation of Theorem D.8 and Theorem*
963 *D.11, the mean squared error of STITCH-OPE is bounded by*

$$\mathbb{E}_{\hat{p}_\pi} \left[(\hat{J} - J(\pi))^2 \right] \leq \left(\frac{2B_w}{1 - \gamma^w} \kappa^w \delta_\beta \right)^2 + 10 \left(\frac{T}{w} \right)^2 B_w^2 \kappa^w \delta_\beta + \frac{8B_w^2}{1 - \gamma^{2w}} \kappa^w \delta_\beta + \text{Var}_{p_\pi}(J).$$

964 *Proof.* We start by adapting the standard bias-variance decomposition to our setting:

$$\begin{aligned}
\mathbb{E}_{\hat{p}_\pi} \left[(\hat{J} - J(\pi))^2 \right] &= \mathbb{E}_{\hat{p}_\pi} \left[(\hat{J} - \mathbb{E}_{\hat{p}_\pi}[\hat{J}] + \mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi))^2 \right] \\
&= \mathbb{E}_{\hat{p}_\pi} \left[(\hat{J} - \mathbb{E}_{\hat{p}_\pi}[\hat{J}])^2 \right] + \mathbb{E}_{\hat{p}_\pi} \left[(\mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi))^2 \right] \\
&\quad + \mathbb{E}_{\hat{p}_\pi} \left[(\hat{J} - \mathbb{E}_{\hat{p}_\pi}[\hat{J}]) (\mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi)) \right] \\
&= \text{Var}_{\hat{p}_\pi}(\hat{J}) + \text{Bias}_{\hat{p}_\pi}(\hat{J})^2 + (\mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi)) (\mathbb{E}_{\hat{p}_\pi}[\hat{J}] - J(\pi)) \\
&= \text{Var}_{\hat{p}_\pi}(\hat{J}) + \text{Bias}_{\hat{p}_\pi}(\hat{J})^2,
\end{aligned}$$

965 since the last term is zero. Plugging in the bounds of Theorems D.8 and D.11 completes the proof. \square

966 E Pseudocode

967 A high-level pseudocode of conditional diffusion model training in STITCH-OPE is provided as
968 Algorithm 1. A pseudocode of the off-policy evaluation subroutine for a single rollout is provided as
969 Algorithm 2. Empirically, we have found that per-term normalization of the guidance function (line
970 9) resulted in more consistent performance, and allowed the guidance coefficients α and λ to be more
971 easily tuned.

972 F Domains

973 We include experiments on the medium datasets from the D4RL offline suite [12], and Pendulum
974 and Acrobot domains from the OpenAI Gym suite [4]. We set the evaluation horizon to $T = 768$ for
975 D4RL, $T = 256$ for Acrobot and $T = 196$ for Pendulum, and we use $\gamma = 0.99$ in all experiments.
976 Furthermore, Acrobot uses a discrete action space and is incompatible with our method, so we
977 modified the domain to take continuous actions. Table 4 summarizes the key properties of each
978 domain.

Algorithm 1 Conditional Diffusion Model Training in STITCH-OPE

Require: diffusion model $\epsilon_\theta(\tau, k|s)$, behavior data \mathcal{D}_β , $w \geq 0$, learning rate $\eta > 0$, $\{\sigma_k\}_{k=1}^K$ and $\{\alpha_k\}_{k=1}^K$ positive

- 1: $\bar{\alpha}_k \leftarrow \prod_{t=1}^k \alpha_t$ **for** $k = 1 \dots K$
- 2: **initialize** θ randomly
- 3: **repeat**
- 4: **sample** length- w sub-trajectory $\tau^0 = (s_0, a_0, s_1, \dots, s_w)$ from \mathcal{D}_β
- 5: **sample** $k \sim \text{Uniform}(\{1, \dots, K\})$ ▷ Sample denoising time step k
- 6: **sample** $\epsilon \sim \mathcal{N}(0, I)$ ▷ Sample pure noise sub-trajectory
- 7: $\nabla_\theta \mathcal{L}(\theta) \leftarrow \nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_k} \tau^0 + \sigma_k \epsilon, k|s_0)\|^2$ ▷ Gradient descent step on θ
- 8: $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta)$
- 9: **until** converged
- 10: **return** ϵ_θ

Algorithm 2 Off-Policy Evaluation in STITCH-OPE

Require: diffusion model $\epsilon_\theta(\tau, k|s)$ (Algorithm 1), empirical reward function $\hat{R}(s, a)$, behavior policy $\beta(a|s)$, target policy $\pi(a|s)$, $\alpha \geq 0$, $\lambda \geq 0$, $w \geq 0$ (divides T), $\{\sigma_k\}_{k=1}^K$ and $\{\alpha_k\}_{k=1}^K$ positive

- 1: $\hat{J} \leftarrow 0$
- 2: **sample** $s_0^0 \sim d_0$ ▷ Sample initial state
- 3: **for** $t = 0$ **to** $T/w - 1$ **do** ▷ Generation for decision epochs wt to $w(t+1)$
- 4: **sample** $\tau_{wt:w(t+1)}^K \sim \mathcal{N}(0, I)$ ▷ Sample pure noise sub-trajectory
- 5: **for** $k = K$ **to** 1 **do** ▷ Denoising step k
- 6: $\mu_{k-1} \leftarrow \frac{1}{\sqrt{\alpha_k}} \left(\tau_{wt:w(t+1)}^k - \frac{1-\alpha_k}{\sigma_k} \epsilon_\theta(\tau_{wt:w(t+1)}^k, k|s_{wt}^0) \right)$ ▷ Mean of diffusion
- 7: $g_k^\pi \leftarrow \sum_{u=wt}^{w(t+1)-1} \nabla_\tau \log \pi(a_u^k | s_u^k)$ ▷ Compute π guidance term
- 8: $g_k^\beta \leftarrow \sum_{u=wt}^{w(t+1)-1} \nabla_\tau \log \beta(a_u^k | s_u^k)$ ▷ Compute β guidance term
- 9: $g_k \leftarrow \alpha(g_k^\pi / \|g_k^\pi\|_2) - \lambda(g_k^\beta / \|g_k^\beta\|_2)$ ▷ Compute normalized guidance
- 10: **sample** $\tau_{wt:w(t+1)}^{k-1} \sim \mathcal{N}(\mu_k + \sigma_k^2 g_k, \sigma_k^2 I)$ ▷ Apply guided diffusion step
- 11: **end for**
- 12: $\hat{J} \leftarrow \hat{J} + \sum_{u=wt}^{w(t+1)-1} \gamma^u \hat{R}(s_u^0, a_u^0)$ ▷ Update π return using denoised $\tau_{wt:w(t+1)}^0$
- 13: **end for**
- 14: **return** \hat{J}

Description	Hopper	Walker	HalfCheetah	Pendulum	Acrobot
state dimension	11	17	17	3	6
action dimension	3	6	6	1	3
range of action	$[-1, 1]$	$[-1, 1]$	$[-1, 1]$	$[-2, 2]$	$[-1, 1]$
rollout length T	768	768	768	196	256
discount factor γ	0.99	0.99	0.99	0.99	0.99

Table 4: Properties of D4RL [12] and OpenAI Gym [4] benchmark problems.

979 G Policies

980 **D4RL Offline Suite** Behavior and target policies and their trained procedures are described in
981 [13], and the policy parameters are borrowed from the official repository at https://github.com/google-research/deep_ope (Apache 2.0 licensed). The 10 target policies of varying ability,
982 $\pi_{\theta_1}, \pi_{\theta_2}, \dots, \pi_{\theta_{10}}$, are obtained by checkpointing the policy parameters $\theta_1, \theta_2 \dots \theta_{10}$ at various points
983 during training. Each target policy network models the action probability distribution $\pi_i(a|s)$ using a
984 set of independent Gaussian distributions, predicting the mean and variance (μ_i, σ_i^2) of each action
985 component a_i independently. This allows the score function of the target policy to be easily computed.
986 As discussed in the main text, all policies are derived from the medium datasets in all experiments.
987

988 **OpenAI Gym** We model target policies $\pi_1, \pi_2 \dots \pi_5$ as MLPs and train them in each environment
 989 following the Twin-Delayed DDPG (TD3) [8] algorithm. The total training time is set to 50000 steps,
 990 and we checkpoint policies every 5000 steps. The behavior policy is set to the target policy π_3 . The
 991 complete list of hyper-parameters is provided in Table 5.

Description	Value
number of hidden layers in actor and critic	2
number of neurons per layer in actor and critic	256
hidden activation function	ReLU
output activation function	tanh
Gaussian noise for exploration	0.1
noise added to target policy during critic update	0.2
target noise clipping	0.5
frequency of delayed policy updates	2
moving average of target θ'	0.005
learning rate of Adam optimizer	0.0003
batch size	256
replay buffer size	1000000

Table 5: Hyper-parameters for training target policies on OpenAI Gym domains.

992 **Bounded Action Space** Since the action spaces for all domains are compact bounded intervals,
 993 we need to restrict the action space of the policy networks during evaluation. We accomplish this by
 994 applying the tanh transformation to each Gaussian action distribution and then scaling the result to
 995 the required range. Note that this transformation constrains the action probability distribution of all
 996 policies to a bounded range, and thus satisfies the requirement of Assumption 3.1.

997 H Baselines

998 The following model-free baseline methods were chosen for empirical comparison with STITCH-
 999 OPE:

1000 **Fitted Q-Evaluation (FQE)** [21] evaluates a target policy π by estimating its Q-value function
 1001 $Q_\theta(s, a)$ using a neural network. The loss function for θ is

$$\mathcal{L}_{FQE}(\theta) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}_\beta, a' \sim \pi(\cdot|s')} \left[(Q_\theta(s, a) - r - \gamma Q_\theta(s', a'))^2 \right].$$

1002 We follow [56, 52] and learn a target Q-network $Q_{\theta'}(s, a)$ in parallel for added stability. We use
 1003 the AdamW algorithm [54] for optimizing the loss function in a minibatched setting, with gradient
 1004 clipping applied to limit the norm of each gradient update to 1. The complete list of hyper-parameters
 1005 used is provided in Table 6.

Description	Hopper	Walker	HalfCheetah	Pendulum	Acrobot
number of hidden layers	2	2	2	2	2
number of neurons per layer	500	500	500	256	100
hidden activation function	sigmoid	sigmoid	sigmoid	sigmoid	sigmoid
learning rate of AdamW optimizer	0.001	0.003	0.00003	0.003	0.001
moving average of target θ'	0.05	0.05	0.001	0.005	0.05
training epochs (passes over data set)	100	50	70	100	200
batch size	512	256	256	128	512

Table 6: Hyper-parameters for Fitted Q-Evaluation (FQE).

1006 **Doubly Robust (DR)** [18, 36] leverages both importance sampling and value function estimation
 1007 to construct a combined estimate that is accurate when either one of the individual estimates is
 1008 correct. First, we define an estimate $\hat{Q}(s, a)$ of the Q-value function of policy π , and let $\hat{V}(s) =$

1009 $\mathbb{E}_{a \sim \pi(\cdot|s)}[\hat{Q}(s, a)]$ be the corresponding value estimate. We also define $\rho_t = \frac{\pi(a_t|s_t)}{\beta(a_t|s_t)}$ as the policy
 1010 ratio at step t . Then, the DR estimator is defined recursively as

$$V_{DR}^{t+1} = \hat{V}(s_t) + \rho_t \left(r_t + \gamma V_{DR}^t - \hat{Q}(s_t, a_t) \right),$$

1011 such that the policy value estimate $\hat{J}_{DR}(\pi) = V_{DR}^0$. We parameterize both $\hat{Q}(s, a)$ and $\hat{V}(s)$ as
 1012 MLPs and train them using AdamW in a mini-batched setting. Similar to FQE, we also update a
 1013 target value network to improve convergence. The full list of hyper-parameters is provided in Table 7.

Description	Hopper	Walker	HalfCheetah	Pendulum	Acrobot
number of hidden layers	2	2	2	2	2
number of neurons per layer	500	500	500	256	100
hidden activation function	sigmoid	sigmoid	sigmoid	sigmoid	sigmoid
learning rate of AdamW optimizer	0.0003	0.003	0.003	0.003	0.00003
moving average of target θ'	0.05	0.05	0.05	0.05	0.001
training epochs (passes over data set)	50	50	50	100	100
batch size	32	256	512	256	128

Table 7: Hyper-parameters for Doubly Robust (DR) estimation.

1014 **Importance Sampling (IS)** [32] evaluates the target policy by importance weighting the full
 1015 trajectory returns in the behavior dataset, i.e.

$$\hat{J}_{IS}(\pi) = \mathbb{E}_{\tau \sim p_\beta} \left[\left(\prod_{t=0}^{T-1} \frac{\pi(a_t|s_t)}{\beta(a_t|s_t)} \right) \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right].$$

1016 It requires access to the target and behavior policy probabilities in order to compute the weighting.
 1017 Specifically, we use the *per-decision* variant of IS (PDIS), i.e.

$$\hat{J}_{PDIS}(\pi) = \mathbb{E}_{\tau \sim p_\beta} \left[\sum_{t=0}^{T-1} \gamma^t \left(\prod_{u=0}^t \frac{\pi(a_u|s_u)}{\beta(a_u|s_u)} \right) R(s_t, a_t) \right],$$

1018 which has lower variance than IS.

1019 **Density Ratio Estimation (DRE)** [30] estimates the ratio $w(s, a) = d^\pi(s, a)/d^\beta(s, a)$ of the
 1020 discounted state-action occupancies of the target policy π relative to the behavior policy β . The
 1021 *discounted state-action occupancy* of a policy $\mu \in \{\beta, \pi\}$ is defined as

$$d^\mu(s, a) = \lim_{T \rightarrow \infty} \frac{\sum_{t=0}^T \gamma^t p(s_t = s, a_t = a | \mu)}{\sum_{t=0}^T \gamma^t},$$

1022 where $p(s_t = s, a_t = a | \mu)$ indicates the probability of sampling state-action pair (s, a) from μ at
 1023 time step t . We also tested the variants of DICE [42] but found their performance to be unsatisfactory,
 1024 so they have been omitted from the study. The target policy value is estimated as

$$\hat{J}(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{(s, a, r) \sim \mathcal{D}_\beta} [w(s, a) \cdot r].$$

1025 $w(s, a)$ is parameterized as a feedforward neural network and its parameters are trained using Adam
 1026 in a mini-batched setting. Fixed hyper-parameters necessary to reproduce the experiment are listed in
 1027 Table 8. Additionally, since the method requires a kernel function to be specified, we use a Gaussian
 1028 kernel $k(x, x') = \exp(-\eta \|x - x'\|^2)$, where x and x' are concatenations of the (standardized) state
 1029 and action vectors. Since this requires setting a kernel bandwidth $\eta > 0$ which affects the overall
 1030 performance significantly, we run this baseline for different values $\eta \in \{0.01, 0.1, 1, 10, 100\}$ and
 1031 report the best performing result (according to log-RMSE).

1032 The following model-based baseline methods were also chosen for empirical comparison with
 1033 STITCH-OPE. They were chosen to determine the benefits of STITCH-OPE compared to fully
 1034 autoregressive sampling, i.e. $w = 1$, and non-autoregressive sampling, i.e. $w = T$.

Description	Value
number of hidden layers of $w(s, a)$	2
number of neurons per layer of $w(s, a)$	256
hidden activation function	Leaky ReLU
output activation function	SoftPlus
learning rate of Adam optimizer	0.001
training epochs (passes over the data set)	20 (D4RL), 200 (Gym)
batch size	512

Table 8: Hyper-parameters for Density Ratio Estimation (DRE) [30].

Model-Based (MB) [18, 39] consists of learning dynamics $\hat{P}(s'|s, a)$, reward function $\hat{R}(s, a)$ and termination function $\hat{D}(s)$ trained on the behavior dataset to directly approximate the data-generating distribution of the target policy, $p_\pi(\tau)$. \hat{P} directly predicts the next state s' given the current state s and action a . Both \hat{P} and \hat{R} can be found by solving a standard nonlinear regression problem, and \hat{D} can be found by solving a binary classification problem trained on termination flags in the behavior dataset. We parameterize all functions as nonlinear MLPs and obtain their optimal parameters using Adam in a mini-batched setting. Once we obtain their optimal parameters, we estimate the target policy return by generating 50 length- T rollouts from the estimated model, and average their empirical cumulative returns. The necessary hyper-parameters are described in Table 9.

Description	Value
number of hidden layers	3
number of neurons per layer	500
hidden activation function	ReLU
learning rate of Adam optimizer	0.0003
training epochs (passes over data set)	100
batch size	1024

Table 9: Hyper-parameters for Model-Based (MB) estimation.

Policy-Guided Diffusion (PGD) [15] takes a generative approach by simulating target policy trajectories using a guided diffusion model. We follow the original implementation by training a diffusion model on the behavior data, using the official implementation located at <https://github.com/EmptyJackson/policy-guided-diffusion> (MIT licensed). We then generate 50 full-length trajectories from the model using guided diffusion [17] with the guidance function $g_{simple}(\tau) = \nabla_\tau \sum_t \log \pi(a_t|s_t)$, using which we estimate the empirical return of the target policy. All hyper-parameters for training the diffusion models are fixed as per the original paper and codebase (see Appendix A therein for details). However, we found that the policy guidance coefficient α and guidance normalization both have significant effects on performance, thus we ran PGD for different choices of $\alpha \in \{0.001, 0.01, 0.1, 1.0, 10, 100, 1000\}$ with and without guidance normalization, and report the best performing result (according to log-RMSE).

I Metrics

Let π_1, \dots, π_{10} be the target policies, $\hat{J}_1(\pi_i), \hat{J}_2(\pi_i), \dots, \hat{J}_5(\pi_i)$ be the estimates of the target policy values across the 5 seeds, and $J(\pi_1), \dots, J(\pi_{10})$ be the target policy values estimated using 300 rollouts collected by running the target policies in the environments.

The following metrics were used to quantify and compare the performance of STITCH-OPE and all metrics:

Log Root Mean Squared Error (LogRMSE) This is defined as the log root mean squared error using the estimates $\hat{J}_j(\pi_1), \dots, \hat{J}_j(\pi_{10})$ and the ground truth returns $J(\pi_1), \dots, J(\pi_{10})$, averaged

1063 across seeds $j = 1 \dots 5$. Mathematically,

$$\frac{1}{5} \sum_{j=1}^5 \log \sqrt{\frac{1}{10} \sum_{i=1}^{10} (\hat{J}_j(\pi_i) - J(\pi_i))^2}.$$

1064 **Spearman (Rank) Correlation** This is defined as the Spearman correlation [61] between the
 1065 estimates $\hat{J}_j(\pi_1), \dots, \hat{J}_j(\pi_{10})$ and the ground truth returns $J(\pi_1), \dots, J(\pi_{10})$, averaged across
 1066 seeds $j = 1 \dots 5$.

1067 **Regret@1** This is defined as the absolute difference in return between the best policy selected using
 1068 the baseline policy returns $\hat{J}_j(\pi_i)$ and the policy selected according to the ground truth estimates
 1069 $J(\pi_i)$, averaged across seeds $j = 1 \dots 5$, i.e:

$$\frac{1}{5} \sum_{j=1}^5 \left| J(\pi_{i_j^{max}}) - \max_{i=1 \dots 10} J(\pi_i) \right|, \quad \text{where } i_j^{max} = \operatorname{argmax}_{i=1 \dots 10} \hat{J}_j(\pi_i).$$

1070 **Normalization** In order to compare metrics consistently across environments, we follow [13] and
 1071 use the normalized policy values:

$$\frac{\hat{J}_j(\pi_i) - V_{min}}{V_{max} - V_{min}}, \quad \text{where } V_{min} = \min_i J(\pi_i), \quad V_{max} = \max_i J(\pi_i),$$

1072 where V_{min} and V_{max} are the minimum and maximum target policy values, respectively.

1073 **Error Bars** All tables and figures report error bars defined as \pm one standard error, i.e. $\hat{\sigma}/\sqrt{n}$
 1074 where $\hat{\sigma}$ is the empirical standard deviation of each metric value across seeds and n is the number of
 1075 seeds (fixed to 5 for all experiments).

1076 J STITCH-OPE Training and Hyper-Parameter Details

1077 We follow the configuration used in [17] for training the diffusion model, including architecture,
 1078 optimizer, and noise schedule. Specifically, we parameterize the diffusion process ϵ as a UNet
 1079 architecture with residual connections [60], trained with a cosine learning rate schedule [55]. The list
 1080 of training hyper-parameters is provided in Table 10. The reward predictor $\hat{R}(s, a)$ is a two-layer
 1081 MLP with ReLU activations and 32 neurons per hidden layer, and is trained using Adam with a
 1082 learning rate of 0.001 and batch size of 64.

Description	Value
diffusion architecture	UNet
learning rate of Adam optimizer	0.0003
training epochs (passes over the data set)	150
batch size	128
training steps per epoch	5000 (D4RL), 2000 (Gym)
guidance coefficient for π , i.e. α	0.5 (D4RL), 0.1 (Gym)
guidance coefficient ratio for β , i.e. $\frac{\lambda}{\alpha}$	0.5 (D4RL), 1 (Gym)
window size of sub-trajectories, i.e. w	8 (D4RL), 16 (Gym)

Table 10: Hyper-parameters for STITCH-OPE.

1083 **Guidance Coefficients** For Gym domains, we use $\alpha = \lambda = 1$, corresponding to the theoretically
 1084 justified guidance function in Equation 8, assuming low distribution shift. For D4RL tasks, we use
 1085 tempered values $\alpha = 0.5$ and $\lambda = 0.25$ to improve sample stability and regularization, which we
 1086 found empirically helpful in higher-dimensional settings.

1087 **Sub-Trajectory Length** We use $w = 16$ for Gym domains and $w = 8$ for most D4RL tasks. For
 1088 HalfCheetah, we reduce to $w = 4$ due to the environment’s fast dynamics, which caused degradation
 1089 in stitching fidelity with longer sub-trajectories. Also to add stability for cheetah, we set the clip
 1090 denoised flag to True during the backward diffusion process.

1091 K Diffusion Policy Training and Evaluation

1092 We follow [40] and parameterize each target policy π'_i , $i = 1 \dots 10$ as a conditional diffusion model
 1093 $\epsilon_{\phi_i}(a^k, k|s)$, whose parameters ϕ_i are learned by optimizing the behavior cloning objective (compare
 1094 with (1))

$$\mathcal{L}(\phi_i) = \mathbb{E}_{k, \epsilon \sim \mathcal{N}(0, I), s \sim \mathcal{D}_\beta, a \sim \pi_i(\cdot|s)} [\|\epsilon - \epsilon_{\phi_i}(a^k, k|s)\|^2].$$

1095 In order to use the fine-tuned $\epsilon_{\phi_i}(a^k, k|s)$ as a guidance function for off-policy evaluation in STITCH-
 1096 OPE, we use the following equivalence between score-based models and denoising diffusion [9]
 1097 (extended trivially to the conditional setting)

$$\nabla_a \log \pi'_i(a|s)|_{a=a^k} = -\frac{\epsilon_{\phi_i}(a^k, k|s)}{\sigma_k}.$$

1098 Specifically, this expression cannot be calculated at $k = 0$ since $\sigma_0 = 0$ using the standard param-
 1099 eterization of diffusion models, so we approximate it at $k = 1$ and use the resulting gradient in
 1100 STITCH-OPE.

1101 We implement the diffusion model using the CleanDiffuser package [49] with official repository
 1102 at <https://github.com/CleanDiffuserTeam/CleanDiffuser> (Apache 2.0 licensed). To train
 1103 the diffusion policies, we first generate rollouts from each of the pre-trained target policies in D4RL
 1104 [13], and then minimize the behavior cloning objective $\mathcal{L}(\phi_i)$ above to obtain the diffusion policy
 1105 parameters. The list of relevant hyper-parameters is provided in Table 11.

Description	Value
embedding dimension	64
hidden layer dimension	256
learning rate	0.0003
diffusion time steps	32
EMA rate	0.9999
total training steps	10000
number of transitions to generate for each dataset	1000000
training batch size	256

Table 11: Hyper-parameters for training diffusion policies.

1106 L Additional Experiments

1107 L.1 Sensitivity to Guidance Coefficients α and λ

1108 We evaluate STITCH-OPE across different choices of the guidance coefficients α and λ , and plot
 1109 the resulting trends in Figure 5 for Hopper and Figure 6 for Walker2D. Each plot is generated by
 1110 applying bicubic interpolation to the grid evaluations of the Spearman correlation and LogRMSE.
 1111 The optimal coefficient values of α and λ remain consistent across environments. The optimal
 1112 balance for off-policy evaluation is attained by assigning a moderate coefficient for the target policy
 1113 score α (i.e. $\alpha < 1$) and a smaller but positive coefficient to the behavior policy score λ , i.e.
 1114 $0 < \lambda < \alpha$. Recall that λ controls the amount of distribution shift we are willing to accept during
 1115 guided trajectory generation. $\lambda = 1$ is theoretically unbiased, but potentially under-regularized and
 1116 leads to dynamically infeasible (high-variance) samples. Meanwhile, $\lambda = 0$ is often over-regularized
 1117 and leads to trajectories that are heavily biased towards the behavior policy $p_\beta(\tau)$. From the plots,
 1118 we see that a moderate amount of regularization is optimal (around 25% of the value of α), which is
 1119 consistent with regularization in supervised machine learning (i.e., regression).

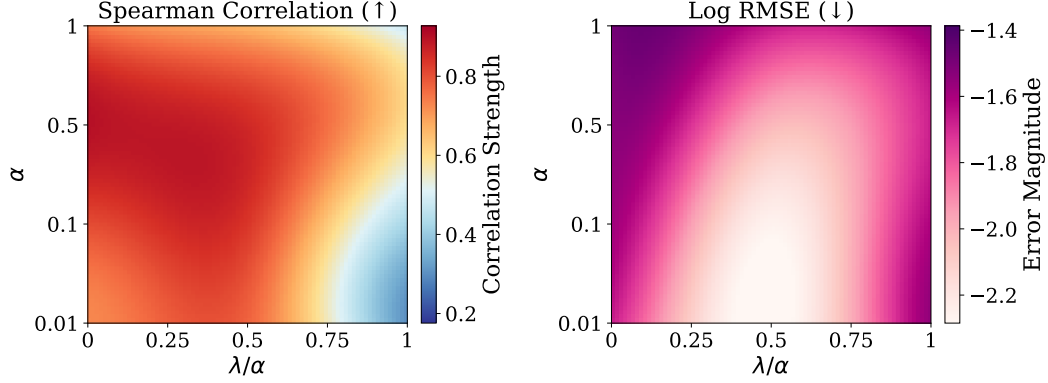


Figure 5: Smoothed performance landscape for Hopper. **Left:** Spearman correlation is largest around $\alpha \in [0.1, 0.5]$, $\lambda \leq 0.5\alpha$. **Right:** The LogRMSE is smallest around $\alpha \in [0.01, 0.5]$, $\lambda \in [0.25\alpha, 0.75\alpha]$. These results confirm the optimal range of λ is $0 < \lambda < \alpha$.

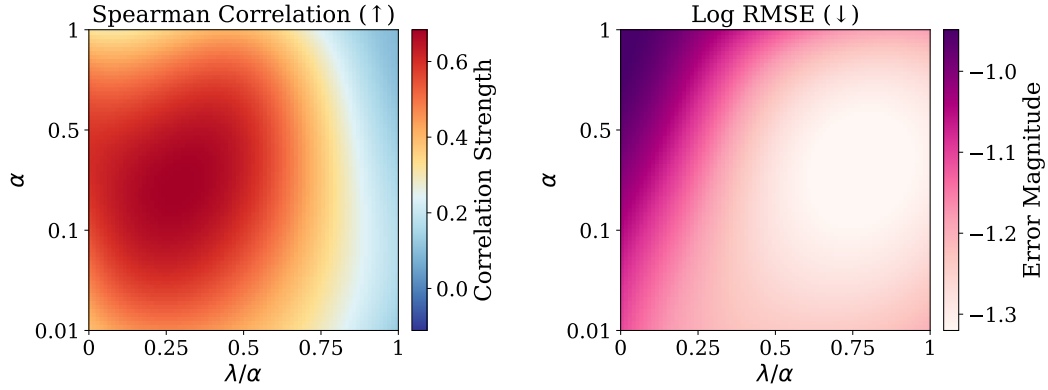


Figure 6: Smoothed performance landscape for Walker2d. Results are generally consistent with Hopper. **Left:** Spearman correlation is largest around $\alpha \in [0.1, 0.5]$, $\lambda \approx 0.25\alpha$. **Right:** The LogRMSE is smallest around $\alpha \in [0.1, 0.5]$, $\lambda \approx 0.75\alpha$. These results confirm the optimal range of λ is $0 < \lambda < \alpha$.

1120 L.2 Sensitivity to Window Size w

1121 To further analyze the sensitivity to w , we evaluate STITCH-OPE across different intermediate values
 1122 of w , and compare the performance according to LogRMSE and Spearman correlation metrics. As
 1123 illustrated in Figure 7 for Hopper and 8 for Pendulum, the best performance is consistently achieved
 1124 using moderate values of w , i.e. $w = 8$ for Hopper and $w = 16$ for Pendulum. As hypothesized in
 1125 the main text, based on our analysis in Section 3.3 and Section 3.4, low values of w provide more
 1126 flexibility when stitching trajectories and thus promote compositionality, but are more susceptible to
 1127 the compounding of errors. High values of w are less susceptible to error compounding but at the
 1128 expense of compositionality and thus less adaptability to distribution shift. In the current ablation
 1129 experiment, it is clear that the best balance between compositionality and error compounding occurs
 1130 using moderate values of w , and the greatest deterioration in performance occurs for very small or
 1131 very large values. It is also important to note that increasing w reduces inference speed due to longer
 1132 trajectory generations per diffusion step, highlighting a practical trade-off between computational
 1133 cost and evaluation accuracy.

1134 L.3 Trajectory Visualizations

1135 We visualize and compare trajectories generated by the guided and unguided versions of STITCH-
 1136 OPE and Policy-Guided Diffusion (PGD) [15] against both random and optimal policies. These
 1137 visualizations highlight differences in the quality of generated trajectories, alignment with target

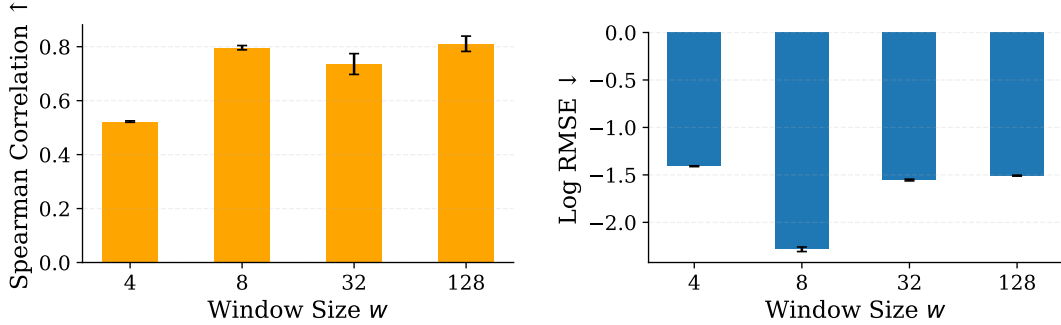


Figure 7: Sensitivity of STITCH-OPE to window size w in the Hopper-v2 environment. **Left:** Spearman rank correlation. **Right:** Log RMSE. Error bars denote one standard error over five random seeds. The overall best performance is attained for $w = 8$, suggesting a good balance between compositionality and error compounding.

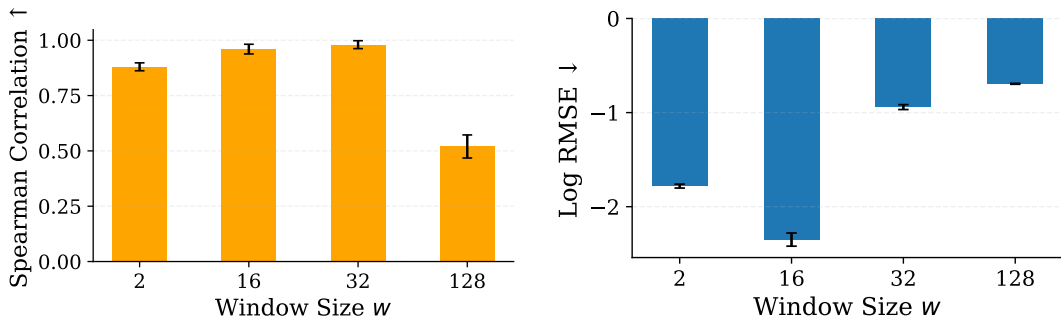


Figure 8: Sensitivity of STITCH-OPE to window size w in the Pendulum-v1 environment. **Left:** Spearman rank correlation. **Right:** Log RMSE. Error bars denote one standard error over five random seeds. The overall best performance is attained for $w = 16$, suggesting a good balance between compositionality and error compounding.

1138 policies, and generalization capabilities across various environments. As shown in Figures 9 and 11,
 1139 STITCH-OPE closely mimics the target policy behavior. On the other hand, PGD performs poorly,
 1140 significantly overestimating the performance of the random policy. Figure 10 further demonstrates
 1141 that STITCH-OPE maintains consistent and robust behavior across policy settings.

1142 M Computing Resources

1143 **Hardware and Software** All experiments were conducted on a local workstation running Ubuntu
 1144 20.04 LTS and Python 3.9, with the following hardware:

- 1145 • 2× NVIDIA RTX 3090 GPUs (24 GB each)
- 1146 • Intel(R) Core(TM) i9-9820X CPU @ 3.30GHz (10 cores / 20 threads)
- 1147 • 128 GB RAM.

1148 **Runtime** Each full training of a diffusion model for a D4RL task took approximately 20 hours to
 1149 complete, depending on environment complexity and rollout length. Each OpenAI Gym task took
 1150 approximately 5 hours. Each evaluation for a D4RL environment took around 18 hours in total (across
 1151 all 5 seeds) to complete, and each OpenAI Gym environment took around 6 hours to complete.

1152 N Related Work

1153 Off-policy evaluation plays a critical role in offline reinforcement learning, enabling the evaluation
 1154 of policies without directly interacting with the environment. OPE has been studied across a wide

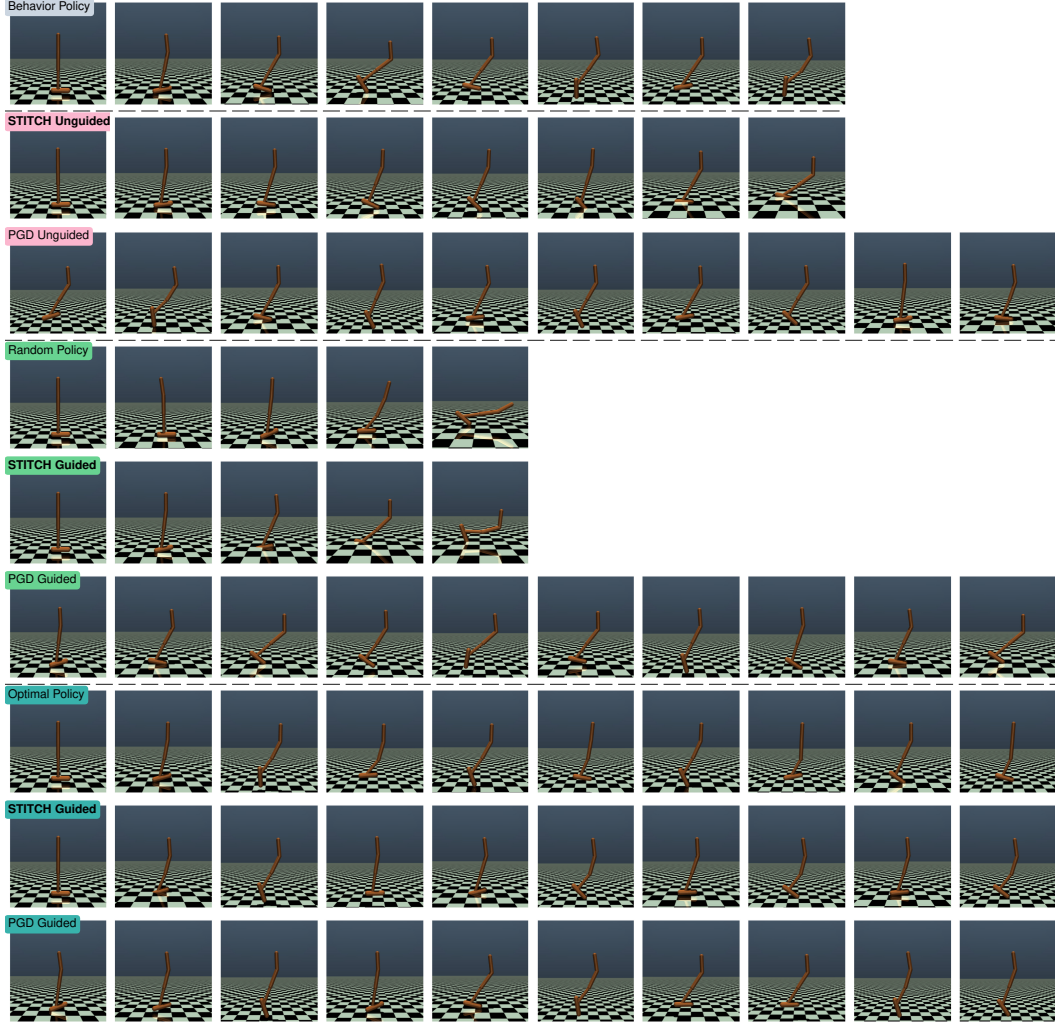


Figure 9: Trajectory visualizations in the Hopper environment. Both STITCH-OPE and PGD track the optimal policy. PGD significantly overestimates the performance of the random policy, while STITCH-OPE correctly models both the state trajectory and the termination.

range of different domains including robotics [51], healthcare [57, 59, 58] and recommender systems [50, 62]. Relevant work includes model-free and model-based OPE approaches, including recent generative methods in offline RL.

Model-Free Methods Model-free methods, such as Importance Sampling (IS) and per-decision Importance Sampling (PDIS) [32] reweight trajectories (or single-step transitions) from the behavior policy to approximate returns under a target policy. However, this class of methods suffers from the so-called “curse of horizon”, in which the variance grows exponentially in the length of the trajectory [24, 26]. Doubly Robust (DR) methods [18, 36, 10] further combine estimation of value functions with importance weights, reducing the overall variance. Distribution-correction methods (DICE) [31, 42, 46] and their variants [24, 30] try to mitigate the curse-of-horizon by performing importance sampling from the stationary distribution of the underlying MDP. However, these methods perform relatively poorly on high-dimensional long-horizon tasks [13].

Model-Based Methods Model-based OPE methods estimate the target policy value by learning approximate transition and reward models from offline data and simulating trajectories under the target policy [18, 20]. These methods have shown strong empirical performance, especially in



Figure 10: Trajectory visualizations in the HalfCheetah environment. STITCH-OPE and PGD both demonstrate consistent behavior across all policy types, highlighting their robust generalization on this task.

1170 continuous control domains [37, 45], but they often suffer from compounding errors during rollouts,
 1171 which can lead to biased estimates in high-dimensional or long-horizon settings [19, 16].

1172 **Offline Diffusion** Inspired by the recent performance of diffusion models across many areas of
 1173 machine learning [14, 9], a new stream of reinforcement learning has emerged which leverages
 1174 diffusion models trained on behavior data [28, 47]. [17, 1] train diffusion models on behavior data
 1175 that can be guided to achieve new goals. [15, 34] apply guided diffusion to offline policy optimization
 1176 by setting the guidance function to be the score of the learned policy, while [63] applies guided
 1177 diffusion to satisfy added safety constraints. Unlike STITCH-OPE, these works do not use negative
 1178 guidance nor stitching, which we found leads to unstable policy values when applied directly for
 1179 offline policy evaluation over a long-horizon. [29] applies DICE to estimate the stationary distribution
 1180 of the underlying MDP, which is used as a guidance function to correct the policy distribution shift for
 1181 offline policy optimization. Unlike STITCH-OPE, this work is not directly applicable to offline policy
 1182 evaluation. Finally, [53] introduces a variant of trajectory stitching for augmenting behavior data, but
 1183 does not apply it for offline policy evaluation. To the best of our knowledge, STITCH-OPE is the first
 1184 work to apply diffusion models to evaluate policies on offline data.



Figure 11: Trajectory visualizations in the Walker2d environment. STITCH-OPE effectively imitates both random and optimal policies. As for the Hopper environment, PGD struggles to correctly imitate the random policy, significantly overestimating its performance.

References

- [48] C. M. Bishop and N. M. Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [49] Z. Dong, Y. Yuan, J. Hao, F. Ni, Y. Ma, P. Li, and Y. Zheng. Cleandiffuser: An easy-to-use modularized library for diffusion models in decision making. *arXiv preprint arXiv:2406.09509*, 2024. URL <https://arxiv.org/abs/2406.09509>.
- [50] M. Dudik, D. Erhan, J. Langford, and L. Li. Doubly robust policy evaluation and optimization. *Statistical Science*, 29(4):485–511, 2014.
- [51] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*, pages 651–673. PMLR, 2018.
- [52] I. Kostrikov and O. Nachum. Statistical bootstrapping for uncertainty estimation in off-policy evaluation. *arXiv preprint arXiv:2007.13609*, 2020.

- 1198 [53] G. Li, Y. Shan, Z. Zhu, T. Long, and W. Zhang. Diffstitch: Boosting offline reinforcement
1199 learning with diffusion-based trajectory stitching. In *International Conference on Machine*
1200 *Learning*, pages 28597–28609. PMLR, 2024.
- 1201 [54] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Con-*
1202 *ference on Learning Representations*, 2019. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Bkg6RiCqY7)
1203 [Bkg6RiCqY7](https://openreview.net/forum?id=Bkg6RiCqY7).
- 1204 [55] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *Interna-*
1205 *tional Conference on Learning Representations*, 2022.
- 1206 [56] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
1207 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-
1208 forcement learning. *nature*, 518(7540):529–533, 2015.
- 1209 [57] S. A. Murphy, M. J. van der Laan, J. M. Robins, and C. P. P. R. Group. Marginal mean models
1210 for dynamic regimes. *Journal of the American Statistical Association*, 96(456):1410–1423,
1211 2001. ISSN 01621459. URL <http://www.jstor.org/stable/3085909>.
- 1212 [58] X. Nie, E. Brunskill, and S. Wager. Learning when-to-treat policies. *Journal of the American*
1213 *Statistical Association*, 116(533):392–409, 2021.
- 1214 [59] A. Raghu, O. Gottesman, Y. Liu, M. Komorowski, A. Faisal, F. Doshi-Velez, and E. Brunskill.
1215 Behaviour policy estimation in off-policy policy evaluation: Calibration matters. *International*
1216 *Conference on Machine Learning: workshop on Causal Machine Learning*, 2018.
- 1217 [60] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image
1218 segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015:*
1219 *18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*,
1220 pages 234–241. Springer, 2015.
- 1221 [61] C. Spearman. The proof and measurement of association between two things. *The American*
1222 *Journal of Psychology*, 15(1):72–101, 1904.
- 1223 [62] G. Theodoropoulos, P. S. Thomas, and M. Ghavamzadeh. Personalized ad recommendation systems
1224 for life-time value optimization with guarantees. In *International Joint Conference on Artificial*
1225 *Intelligence*, 2015. URL <https://api.semanticscholar.org/CorpusID:8081523>.
- 1226 [63] Y. Zheng, J. Li, D. Yu, Y. Yang, S. E. Li, X. Zhan, and J. Liu. Safe offline reinforcement learning
1227 with feasibility-guided diffusion model. In *The Twelfth International Conference on Learning*
1228 *Representations*, 2024. URL <https://openreview.net/forum?id=j5JvZCaDMO>.