

SAMPLING FROM ENERGY-BASED POLICIES USING DIFFUSION

Anonymous authors

Paper under double-blind review

ABSTRACT

Energy-based policies offer a flexible framework for modeling complex, multimodal behaviors in reinforcement learning (RL). In maximum entropy RL, the optimal policy is a Boltzmann distribution derived from the soft Q-function, but direct sampling from this distribution in continuous action spaces is computationally intractable. As a result, existing methods typically use simpler parametric distributions, like Gaussians, for policy representation — limiting their ability to capture the full complexity of multimodal action distributions. In this paper, we introduce a diffusion-based approach for sampling from energy-based policies, where the negative Q-function defines the energy function. Based on this approach, we propose an actor-critic method called Diffusion Q-Sampling (DQS) that enables more expressive policy representations, allowing stable learning in diverse environments. We show that our approach enhances exploration and captures multimodal behavior in continuous control tasks, addressing key limitations of existing methods.

1 INTRODUCTION

Deep reinforcement learning (RL) is a powerful paradigm for learning complex behaviors in diverse domains, from strategy-oriented games (Silver et al., 2016; Berner et al., 2019; Schrittwieser et al., 2020) to fine-grained control in robotics (Kober et al., 2013; Sünderhauf et al., 2018; Wu et al., 2023). In the RL framework, an agent learns to make decisions by interacting with an environment and receiving feedback in the form of reward. The agent aims to learn a policy that maximizes the cumulative sum of rewards over time by exploring actions and exploiting known information about the environment’s dynamics.

The parameterization of the policy is a crucial design choice for any RL algorithm. Under the conventional notion of optimality, under full observability, there always exists an optimal deterministic policy that maximizes the long-term (discounted) return (Sutton & Barto, 2018). However, this is only true when the agent has explored sufficiently and has nothing to learn about the environment. Exploration requires a stochastic policy, to experiment with different potentially rewarding actions. Moreover, even in the exploitation phase, there may be more than one way of performing a task and we might be interested in mastering all of them. This diversification is motivated by the robustness of the resulting stochastic policy to environment changes; if certain pathways for achieving a task become infeasible due to a change of the dynamics or reward, some may remain feasible, and the agent has an easier time in adapting to this change by exploiting and improving the viable options. This argument also suggests that such policies can also serve as effective initialization for fine-tuning on specific tasks.

While exploration, diversity and robustness motivate stochastic policies, representing such policies in continuous action spaces remains challenging. As a result, stochasticity is often introduced by noise injection (Lillicrap et al., 2015) or using an arbitrary parametric family (Schulman et al., 2015) which lacks expressivity.

Orthogonal to the difficulty of representing such policies is their training objective; policies are often optimized to maximize the Q-function, and stochasticity is introduced to encourage exploration as an afterthought. However, our argument for stochasticity favours multi-modal policies; instead of learning the *single best* way to solve a task, we want to learn *all reasonably good* ways to solve the task.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

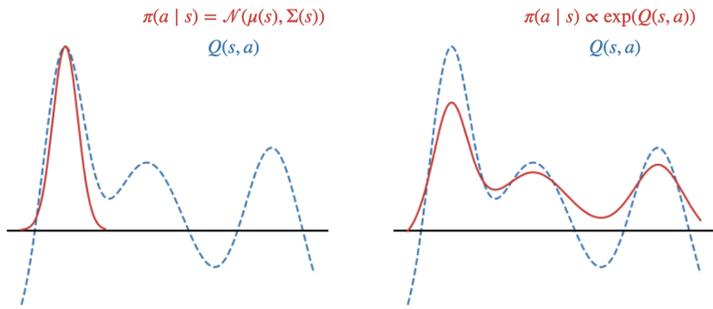


Figure 1: Illustration comparing Gaussian and Boltzmann policies.

We address both of these issues using diffusion-based sampling from a Boltzmann policy. Diffusion models offer a potential solution to the policy parameterization problem since they are expressive and can produce high-quality samples from complex distributions. Indeed, they have been extensively applied to solve sequential decision-making tasks, especially in offline settings where they can model multimodal datasets from suboptimal policies or diverse human demonstrations. A few studies have applied these models in the online setting, focusing on deriving training objectives for policy optimization via diffusion. Yang et al. (2023); Psenka et al. (2023) use the gradient of the Q-function to guide action sampling, however, the exact form of the policy is unspecified and it is unknown what distribution the diffusion models sample from.

We consider the problem of explicitly sampling from energy-based policies of the form,

$$\pi(a | s) \propto \exp(Q^\pi(s, a)).$$

This is also known as the Boltzmann distribution of the Q-function. The optimal policy in the maximum entropy RL framework is also known to be of this form, except it uses the soft Q-function (Haarnoja et al., 2017). Such a policy has several benefits. First, it offers a principled way to balance exploration and exploitation for continuous action spaces. By sampling from this distribution, the policy still prioritizes actions with high Q-values but also has a non-zero probability of sampling sub-optimal actions. While the use of Boltzmann policies is common in the discrete setting, it is challenging in continuous spaces; this sampling problem is often tackled with Markov Chain Monte Carlo (MCMC) techniques, which can be computationally expensive and suffer from exponential mixing time. Second, this formulation naturally incorporates multimodal behavior, since the policy can sample one of multiple viable actions at any given state.

Our contributions in this work are as follows:

- We propose a novel framework for sequential decision-making using diffusion models for sampling from energy-based policies.
- We propose an actor-critic algorithm for training diffusion policies based on our framework.
- We demonstrate the effectiveness of this approach for learning multimodal behaviors in maze navigation tasks.
- We demonstrate improved sample efficiency in continuous control tasks owing to better exploration capabilities of our method.

2 RELATED WORK

Our work is related to two distinct sub-areas of reinforcement learning: the relatively new and actively explored line of work on applying diffusion models in the RL setting, and the classical maximum entropy RL framework.

Diffusion models in RL. Early work on applying diffusion models for RL was focused on behavior cloning in the offline setting (Chi et al., 2023). This setting more closely matches the original purpose of diffusion models - to match a distribution to a given dataset. Janner et al. (2022); Ajay et al. use a diffusion model trained on the offline data for trajectory optimization, while Reuss et al.

(2023); Jain & Ravanbakhsh (2024) apply diffusion models to offline goal-reaching tasks by additionally conditioning the score function on the goal. Within the behavior cloning setting, there is some existing work on learning a stochastic state dynamics model using diffusion (Li et al., 2022).

Beyond behavior cloning, offline RL methods incorporate elements from Q-learning to learn a value function from the offline dataset and leverage the learned Q-function to improve the diffusion policy. A large body of work exists in this sub-field, where the most common approach is to parameterize the policy using a diffusion model and propose different training objectives to train the diffusion policy. Wang et al.; Kang et al. (2024) add a Q-function maximizing term to the diffusion training objective, and Hansen-Estruch et al. (2023) use an actor-critic framework based on a diffusion policy and Implicit Q-learning (Kostrikov et al., 2021). Lu et al. (2023) take an energy-guidance approach, where they frame the problem as using the Q-function to guide the behavior cloning policy to high reward trajectories.

The application of diffusion models has been relatively less explored in the online setting (Ding & Jin, 2023). Yang et al. (2023) modifies actions in the replay buffer based on the gradient of the Q-function, then trains a diffusion model on the modified actions. Psenka et al. (2023) directly trains a neural network to match the gradient of the Q-function, then uses it for sampling. The policy resulting from these methods does not have an explicit form, and it is unknown what exact distribution these diffusion models sample from. In contrast, our method explicitly models the policy as a Boltzmann distribution and samples actions from this distribution.

Maximum entropy RL. In contrast to standard RL, where the goal is to maximize expected returns, in the maximum entropy RL framework, the value function is augmented by Shannon entropy of the policy. Ziebart et al. (2008) applied such an approach in the context of inverse reinforcement learning and Haarnoja et al. (2017) generalized this approach by presenting soft Q-learning to learn energy-based policies. A follow-up work, Haarnoja et al. (2018a), presented the well-known soft actor-critic (SAC) algorithm. This line of work proposes to learn a soft value function by adding the entropy of the policy to the reward. The optimal policy within this framework is a Boltzmann distribution, where actions are sampled based on the exponentiated soft Q-values.

A separate but related line of work on generative flow networks (GFlowNets), originally defined in the discrete case (Bengio et al., 2021; 2023), learns a policy that samples terminal states proportional to the Boltzmann density corresponding to some energy function. They have been extended to the continuous setting (Lahlou et al., 2023) and under certain assumptions, they are equivalent to maximum entropy RL (Tiapkin et al., 2024; Deleu et al.). They can effectively sample from the target distribution using off-policy exploration, however, they encounter challenges in credit assignment and exploration efficiency (Malkin et al., 2022; Madan et al., 2023; Rector-Brooks et al., 2023; Shen et al., 2023). Our approach is distinct as we sample the action at each step from the Boltzmann density of the Q-function, instead of the terminal states based on the reward.

3 PRELIMINARIES

3.1 REINFORCEMENT LEARNING

We consider an infinite-horizon Markov Decision Process (MDP) denoted by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, where the state space \mathcal{S} and action space \mathcal{A} are continuous. $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ denotes the transition probability of the next state $s_{t+1} \in \mathcal{S}$ given the current state $s_t \in \mathcal{S}$ and action $\mathbf{a}_t \in \mathcal{A}$. The reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is assumed to be bounded $r(s, \mathbf{a}) \in [r_{\min}, r_{\max}]$. $\gamma \in (0, 1]$ is the discount factor.

A policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ produces an action for every state $s \in \mathcal{S}$. In the standard RL framework, the objective is to learn a policy that maximizes the expected sum of rewards $\sum_t \mathbb{E}_{a_t \sim \pi(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)} [r(s_t, a_t)]$.

Actor-critic is a commonly used framework to learn such policies. It consists of optimizing a policy (the actor) to choose actions that maximize the action value function, also called the Q-function (the critic). The Q-function is defined as the sum of expected future rewards starting from a given

state-action pair, and thereafter following some policy π ,

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{k=t}^{\infty} r(s_k, a_k) \mid s_t = s, a_t = a \right].$$

The optimal policy is defined as the policy that maximizes the sum of rewards along a trajectory,

$$\pi^* = \arg \max_{\pi} \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} r(s_t, a_t) \right]$$

3.2 DIFFUSION MODELS

Denosing Diffusion. Denosing diffusion (Dinh et al., 2016; Ho et al., 2020; Song et al., 2021) refers to a class of generative models which relies on a stochastic process which progressively transforms the target data distribution to a Gaussian distribution. The time-reversal of this diffusion process gives the generative process which can be used to transform noise into samples from the target data distribution.

The forward noising process is a stochastic differential equation,

$$dx_\tau = -\alpha(\tau)x_\tau d\tau + g(\tau)dw_\tau \quad (1)$$

where w_τ denotes Brownian motion. In this paper, we consider the Variance Exploding (VE) SDE where the decay rate, α , is set to $\alpha(\tau) = 0$. This noising process starts with samples from the target density $x_0 \sim p_0$ and progressively adds noise to them over a diffusion time interval $\tau \in [0, 1]$. The marginal probability distribution at time τ is denoted by p_τ and is the convolution of the target density p_0 with a normal distribution with a time-dependent variance, σ_τ^2 . For the VE setting we consider, these marginal distributions are given by,

$$p_\tau(x_\tau) = \int_0^\tau p_0(x_0) \mathcal{N}(x_\tau; x_0, \sigma_\tau^2) dx_0 \quad (2)$$

where the variance is related to the diffusion coefficient, $g(\tau)$ via $\sigma_\tau^2 = \int_0^\tau g(\xi)^2 d\xi$.

The generative process corresponding to the corresponding to 1 is an SDE with Brownian motion \bar{w}_τ , given by:

$$dx_\tau = [-\alpha(\tau)x_\tau - g(\tau)^2 \nabla \log p_\tau(x_\tau)] d\tau + g(\tau) d\bar{w}_\tau \quad (3)$$

Therefore, to be able to generate data, we need to estimate the score of the intermediate distributions, $\nabla \log p_\tau(x_\tau)$.

Iterated Denosing Energy Matching. Recently, Akhound-Sadegh et al. (2024) proposed an algorithm known as iDEM (Iterated Denosing Energy Matching) for sampling from a Boltzmann-type target distribution, $p_0(x) \propto \exp(-\mathcal{E}(x))$, where \mathcal{E} denotes the energy. iDEM is a diffusion-based neural sampler, which estimates the diffusion score, $\nabla \log p_\tau$ using a Monte Carlo estimator. Given the VE diffusion path defined above, iDEM rewrites the score of the marginal densities as,

$$\begin{aligned} \nabla \log p_\tau &= \frac{\int \nabla \exp(-\mathcal{E}(x_0)) \mathcal{N}(x_\tau; x_0, \sigma_\tau^2) dx_0}{\int \exp(-\mathcal{E}(x_0)) \mathcal{N}(x_\tau; x_0, \sigma_\tau^2) dx_0} \\ &= \frac{\mathbb{E}_{\tilde{x} \sim \mathcal{N}(x_\tau, \sigma_\tau^2)} [\nabla \exp(-\mathcal{E}(x))]}{\mathbb{E}_{\tilde{x} \sim \mathcal{N}(x_\tau, \sigma_\tau^2)} [\exp(-\mathcal{E}(x))]} \end{aligned} \quad (4)$$

Rearranging the equation above leads to the K -sample Monte-Carlo estimator of the score, given by,

$$S_k(x_\tau, \tau) = \nabla \log \sum_{i=1}^K \exp(-\mathcal{E}(\tilde{x}^{(i)})), \quad \tilde{x}^{(i)} \sim \mathcal{N}(x_\tau, \sigma_\tau^2) \quad (5)$$

A score-network, f_ϕ is trained to regress to the MC estimator, S_k . The network is trained using a bi-level iterative scheme: (1) in the outer-loop a replay buffer is populated with samples that are generated using the model and (2) in the inner-loop the network is regressed $S_k(x_\tau, \tau)$ where x_τ are noised samples from the replay buffer.

216 4 METHOD

217
218 Our objective is to learn general policies of the form $\pi(\mathbf{a} \mid \mathbf{s}) \propto \exp(-\mathcal{E}(\mathbf{s}, \mathbf{a}))$, where \mathcal{E} represents
219 an energy function which specifies the desirability of state-action pairs. By setting the Q-function,
220 $Q(\mathbf{s}, \mathbf{a})$ as the negative energy, we get what is known as the Boltzmann policy,
221

$$222 \pi(\mathbf{a} \mid \mathbf{s}; T) = \frac{\exp(\frac{1}{T}Q(\mathbf{s}, \mathbf{a}))}{\int_{\mathbf{a}} \exp(\frac{1}{T}Q(\mathbf{s}, \mathbf{a}))d\mathbf{a}}. \quad (6)$$

223
224
225 Choosing such a policy gives us a principled way to balance exploration and exploitation. Specif-
226 ically, by scaling the energy function with a temperature parameter T and annealing it to zero, we
227 get a policy that initially explores to collect more information about the environment and over time
228 exploits the knowledge it has gained.
229

230 4.1 BOLTZMANN POLICY ITERATION

231
232 We propose a policy iteration process that alternates between learning the Q-function of the current
233 policy and optimizing the policy to sample from the Boltzmann distribution of this Q-function.
234 While Boltzmann policies have been used in Q-learning methods for discrete action spaces, our
235 approach introduces a key difference – we rely on the Q-function of the current policy, rather than
236 the optimal Q-function. Since the policy is not maximizing the Q-function but rather following a
237 Boltzmann distribution, and the Q-function corresponds to the current policy, it is not immediately
238 clear that this iterative process of policy evaluation and improvement will converge to a unique
239 fixed point. Indeed, it is conceivable to encounter policy-value pairs that are quite sub-optimal yet
240 consistent, where the policy is Boltzmann in Q , and Q is the value function of that specific policy.
241 The following results show that these updates still constitute policy iteration and there is a unique
242 fixed point. Moreover, such a fixed point is bound to recover the optimal policy.

243 **Policy Evaluation.** In the policy evaluation step, we wish to estimate the value function associated
244 with a policy π . For a fixed policy, the Q-function is computed iteratively, starting from any function
245 $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and repeatedly applying the Bellman backup operator \mathcal{T}^π given by,
246

$$247 \mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t) \triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim \mathcal{P}, \mathbf{a}_{t+1} \sim \pi} [Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})] \quad (7)$$

248 This is the standard Bellman backup and the usual convergence results for policy evaluation can be
249 applied (Sutton & Barto, 2018).
250

251
252 **Policy Improvement.** In the policy improvement step, we update the policy to match the Boltz-
253 mann density of the current Q-function. The following lemma shows that this choice of policy
254 update guarantees that the new policy is better than the old policy in terms of the Q-values. Here, Π
255 denotes some set of policies, which corresponds in our case to the set of policies represented by the
256 diffusion process.

257 **Lemma 1.** *Let $\pi_{old} \in \Pi$ and π_{new} be the policy from Equation (6) with respect to $Q^{\pi_{old}}$. Then*
258 $Q^{\pi_{new}}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t)$ *for all $\mathbf{s}_t, \mathbf{a}_t \in \mathcal{S}, \mathcal{A}$.*

259
260 *Proof.* See Appendix A.1 □

261
262
263 **Policy Iteration.** By alternating between policy evaluation and policy improvement, we can prove
264 that it converges to the optimal policy within Π .

265 **Theorem 1.** *Alternating between policy evaluation and policy improvement repeatedly from any*
266 $\pi \in \Pi$ *converges to a policy π^* such that $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ for all $\pi \in \Pi$ and $(\mathbf{s}_t, \mathbf{a}_t) \in$*
267 $\mathcal{S} \times \mathcal{S}$.

268
269 *Proof.* See Appendix A.2 □

Algorithm 1: Diffusion Q-Sampling (DQS)

```

270 Initialize: Initialize Q-function parameters  $\theta$ , policy parameters  $\phi$ , target network  $\bar{\theta} \leftarrow \theta$ ,
271 replay buffer  $\mathcal{D}$ 
272
273 for each iteration do
274   // Environment Interaction
275   for each environment step do
276     Observe state  $\mathbf{s}_t$  and sample action  $\mathbf{a}_t$  via reverse diffusion using  $f_\phi$ 
277     Execute  $\mathbf{a}_t$ , observe reward  $r_t$  and next state  $\mathbf{s}_{t+1}$ 
278     Store transition  $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  in  $\mathcal{D}$ 
279   end
280   // Parameter Updates
281   for each gradient step do
282     Sample minibatch  $\mathcal{B} = \{(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})\}$  from  $\mathcal{D}$ 
283     // Update Q-function parameters  $\theta$ 
284     Compute target Q-values:  $\hat{Q}_t = r_t + \gamma Q_{\bar{\theta}}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})$ ,  $\mathbf{a}_{t+1} \sim \pi_\phi(\mathbf{s}_{t+1})$ 
285     Update  $\theta$  by minimizing:  $J(\theta) = \frac{1}{|\mathcal{B}|} \sum_{\mathcal{B}} (Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}_t)^2$ 
286     // Update policy parameters  $\phi$ 
287     for each  $(\mathbf{s}_t, \mathbf{a}_t)$  in  $\mathcal{B}$  do
288       Sample diffusion time  $\tau \sim \mathcal{U}[0, 1]$ 
289       Sample noisy action  $\mathbf{a}_{t,\tau} \sim \mathcal{N}(\mathbf{a}_t, \sigma_\tau^2 \mathbf{I})$ 
290       Sample  $\{\tilde{\mathbf{a}}_t^{(i)}\}_{i=1}^K$ , where  $\tilde{\mathbf{a}}_t^{(i)} \sim \mathcal{N}(\mathbf{a}_{t,\tau}, \sigma_\tau^2 \mathbf{I})$ 
291       Estimate score:  $S_t = \nabla_{\mathbf{a}_{t,\tau}} \log \sum_{i=1}^K \exp(Q_\theta(\mathbf{s}_t, \tilde{\mathbf{a}}_t^{(i)}))$ 
292       Update  $\phi$  by minimizing:  $J(\phi) = \|f_\phi(\mathbf{s}_t, \mathbf{a}_{t,\tau}, \tau) - S_t\|^2$ 
293     end
294     // Update target network
295     Update  $\bar{\theta} \leftarrow \eta \theta + (1 - \eta) \bar{\theta}$ 
296   end
297 end

```

4.2 DIFFUSION Q-SAMPLING

In this section, we propose an off-policy actor-critic algorithm, which we call *Diffusion Q-Sampling* (DQS), based on the above framework. Being an off-policy method means DQS can reuse past interactions with the environment by storing them in a replay buffer \mathcal{D} , improving sample efficiency.

Let Q_θ denote the Q-function and π_ϕ a parametric policy, where θ, ϕ represent the parameters of a neural network. The Q-function is learned using standard temporal difference learning,

$$J(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right], \quad (8)$$

where

$$\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim \mathcal{P}, \mathbf{a}_{t+1} \sim \pi_\phi} [Q_{\bar{\theta}}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})]$$

The target Q-values, \hat{Q} , make use of a target Q-network denoted by $Q_{\bar{\theta}}$, where the parameters $\bar{\theta}$ are usually an exponentially moving average of the Q-network parameters θ . Also, in practice, the expectation over next states \mathbf{s}_{t+1} is estimated using only a single sample.

We parameterize the policy using a diffusion process and use iDEM (Akhound-Sadegh et al., 2024) to sample actions from the target density $\pi(\cdot | \mathbf{s}_t) \propto \exp(Q_\theta(\mathbf{s}_t, \mathbf{a}_t))$.

Forward process. Given $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$, we progressively add Gaussian noise to the action following some noise schedule. Let $\mathbf{a}_{t,\tau}$ denote the noisy action at diffusion step $\tau \in [0, 1]$, such that,

$$\mathbf{a}_{t,0} = \mathbf{a}_t; \quad \mathbf{a}_{t,\tau} \sim \mathcal{N}(\mathbf{a}_t, \sigma_\tau^2 \mathbf{I}).$$

We choose a geometric noise schedule $\sigma_\tau = \sigma_{\min} \left(\frac{\sigma_{\max}}{\sigma_{\min}} \right)^\tau$, where σ_{\min} and σ_{\max} are hyperparameters. We found it sufficient to set $\sigma_{\min} = 10^{-5}$ and $\sigma_{\max} = 1.0$ for all our experiments.

Reverse process. Given noisy action samples, we iteratively denoise them using a learned score function to produce a sample from the target action distribution $\pi(\cdot|\mathbf{s}_t) \propto \exp(Q^\pi(\mathbf{s}_t, \cdot))$. We train a neural network, f_ϕ to match iDEM’s K -sample Monte Carlo estimator of the score, defined in Equation (5), by setting the negative Q-function as the energy function. The score function takes as input the noisy action and diffusion time, while also being conditioned on the current state. The loss function is given by:

$$J(\phi) = \mathbb{E}_{\substack{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}, \tau \sim U[0,1], \\ \mathbf{a}_{t,\tau} \sim \mathcal{N}(\mathbf{a}_t, \sigma_\tau^2 I)}} \left[\left\| f_\phi(\mathbf{s}_t, \mathbf{a}_{t,\tau}, \tau) - \nabla \log \sum_{i=1}^K \exp(Q_\theta(\mathbf{s}_t, \tilde{\mathbf{a}}_t^{(i)})) \right\|^2 \right], \quad (9)$$

where $\tilde{\mathbf{a}}_t^{(i)} \sim \mathcal{N}(\mathbf{a}_{t,\tau}, \sigma_\tau^2 I)$.

Summarizing, to sample an action \mathbf{a}_t for the current state \mathbf{s}_t such that $\pi_\phi(\mathbf{a}_t|\mathbf{s}_t) \propto \exp(Q^{\pi_\phi}(\mathbf{s}_t, \mathbf{a}_t))$, we first sample noise from the prior (corresponding to diffusion time $\tau = 1$) $\mathbf{a}_{t,1} \sim \mathcal{N}(0, \sigma_1^2)$. We then use Equation (3) in the VE setting (i.e. $\alpha(\tau) = 0$) by using the trained score function f_ϕ in place of $\nabla \log p_\tau$ to iteratively denoise samples produce the action sample \mathbf{a}_t . The full algorithm is presented in Algorithm 1.

Temperature. We can incorporate the temperature parameter T from Equation (6) within our framework by simply scaling the Q-function in Equation (9) and regressing to the estimated score of the temperature-scaled Boltzmann distribution. To enable the score network to model this temperature scaling accurately, we additionally condition f_ϕ on the current temperature. Generally, the temperature is set to a high value initially, and is annealed over time such that at $t \rightarrow \infty$, we have $T \rightarrow 0$. In practice, the temperature is annealed to a sufficiently small value for large time steps. This ensures that the policy explores initially and as it collects more information about the environment, starts exploiting more and more as time passes.

5 EXPERIMENTS

We perform experiments to answer the following major questions:

- Can DQS learn multimodal behaviors, i.e., learn multiple ways to solve a task?
- Does DQS offer better exploration-exploitation in continuous control tasks?

Baselines. We test our method against two relevant baseline methods: (1) Soft Actor-Critic (SAC) (Haarnoja et al., 2018a), a maximum entropy RL method that is widely used for continuous state-action spaces, and; (2) Q-Score Matching (QSM) (Psenka et al., 2023), a recent diffusion-based approach that trains a score function to match the gradient of the Q-function and uses this score function to sample actions.

All methods were trained with $250k$ environment interactions and one network update per environment step. For a fair comparison, all policy/score networks are MLPs with two hidden layers of dimension 256 each, and the learning rate for all networks is 3×10^{-4} . We apply the double Q-learning trick in our implementation, a commonly used technique, where two Q-networks are trained independently and their minimum value is used for policy evaluation to avoid overestimation bias.

5.1 GOAL REACHING MAZE

We use a custom maze environment to evaluate the ability of our method to reach multiple goals. The agent is tasked with manipulating a ball to reach some unknown goal position in the maze. The state consists of the ball’s (x, y) position and the velocity vector. The action is the force vector applied to the ball. The initial state of the ball is at the center of the maze, with some noise added for variability. We define two potential goal states for the ball - the top left and the bottom right corners

378
379
380
381
382
383
384
385
386
387
388
389

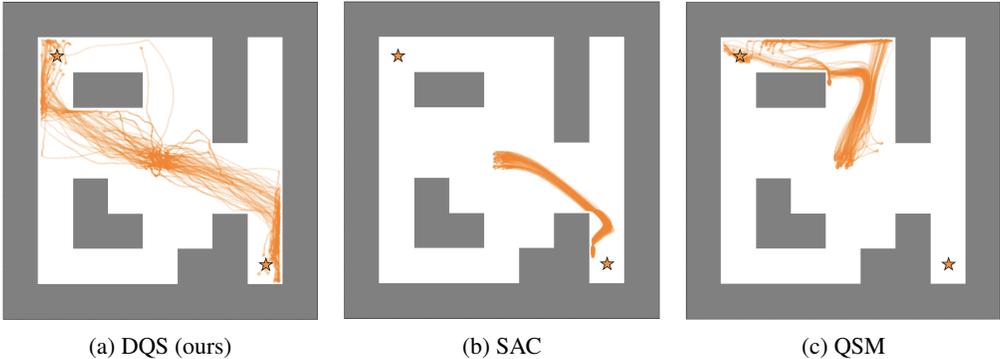


Figure 2: Trajectories for 100 evaluation episodes after 250k training steps.

390
391
392

393 respectively. The negative Euclidean distance between the desired goal and the achieved state gives
394 the reward function.

395
396
397
398
399

For DQS, we used temperature annealing with an initial temperature of $T = 10$, which is decayed
exponentially with the number of training steps to a value of $T = 1$ after 250k steps. SAC uses au-
tomatic temperature tuning (Haarnoja et al., 2018b), where the entropy co-efficient is automatically
tuned using gradient descent to maintain the desired level of entropy.

400
401
402
403
404
405
406

Figure 2 plots the trajectories of the ball over 100 evaluation episodes after 250k training steps. As
seen in Figure 2a, DQS learns to reach both goals, owing to the proposed sampling approach which
can effectively capture multimodal behavior. Moreover, it discovers both paths to reach the top left
goal. In contrast, both SAC (Figure 2b) and QSM (Figure 2c) can only reach one of the goals. Since
SAC models the policy using a Gaussian, there is little variability between different trajectories.
QSM produces slightly more varied behavior, since it also uses a diffusion model to sample actions,
but ultimately fails to learn distinct behaviors.

407
408

5.2 DEEPMIND CONTROL SUITE

409
410
411
412
413
414



415
416

Figure 3: Domains from DeepMind Control Suite considered in our experiments - cheetah, finger,
fish, reacher, hopper, quadruped and walker.

417
418
419
420
421
422

We evaluate the performance of DQS on several continuous control tasks via the DeepMind Control
Suite. We choose eight tasks from different domains to cover tasks of varying complexity and
dynamics. These tasks typically involve controlling the torques applied at the joints of robots to
reach a specific configuration or location, or for locomotion.

423
424
425
426
427
428

Since we are interested in evaluating the ability of DQS to balance exploration and exploitation, we
focus on data efficiency by limiting the number of environment interactions to 250k steps. Figure 4
shows the performance of various methods on these different tasks. On most tasks, DQS performs
on par or outperforms the baseline methods. In particular, on five out of the eight tasks considered
(cheetah-run, finger-spin, fish-swim, reacher-hard and walker-run) DQS reaches higher reward much
faster than competing methods, demonstrating improved exploration.

429
430
431

We use a single fixed temperature of $T = 0.05$ across tasks. Note that as mentioned above, SAC
uses automatic temperature tuning which allows it to influence the policy entropy over the course of
training. The performance of DQS may be further improved by fine-tuning the temperature schedule
on each individual task.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

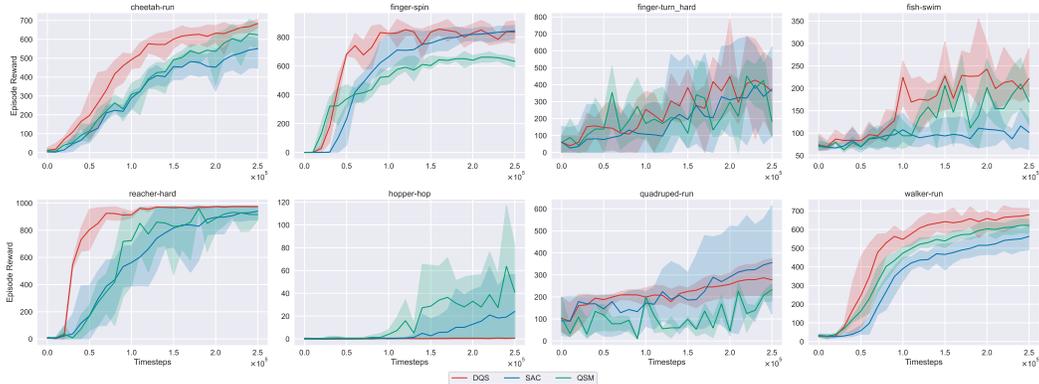


Figure 4: Experimental results on 8 tasks from different domains from the DeepMind Control Suite. Each result is averaged over 10 seeds, with the shaded regions showing minimum and maximum values.

6 DISCUSSION

In this work, we showcase the benefits of using energy-based policies as an expressive class of policies for deep reinforcement learning. Such policies arise in different RL frameworks, but their application has been limited in continuous action spaces owing to the difficulty of sampling in this setting. We alleviate this problem using a diffusion-based sampling algorithm. We show that the process of alternating between updating the Q-function of the current policy and optimizing the policy to sample from the Boltzmann distribution of this Q-function has a unique fixed point which is also the optimal policy. The proposed actor-critic algorithm, Diffusion Q-Sampling (DQS), can sample multimodal behaviors and effectively manage exploration and exploitation.

While diffusion methods offer high expressivity, they often come with increased computation. This is particularly true in the online RL setting, where using a diffusion policy means that each environment step requires multiple function evaluations to sample from the diffusion model. There is a growing body of work on efficient SDE samplers (Jolicœur-Martineau et al., 2021), which aim to reduce the number of function evaluations required to obtain diffusion-based samples while maintaining high accuracy. Incorporating such techniques with Boltzmann policies can greatly reduce the computational cost, especially in high-dimensional state-action spaces.

A crucial aspect of energy-based policies is the temperature parameter, which defines the shape of the sampling distribution. Our method enables annealing of the temperature from some starting value to lower values, as is typically done when applying Boltzmann policies in deep RL. However, this temperature schedule has to be manually tuned. Haarnoja et al. (2018b) proposes an automatic temperature tuning method for SAC, which maintains the temperature so that the entropy of the current policy is close to some target entropy. While such an approach could be applied to DQS in principle, it is computationally expensive to compute the likelihoods of samples under a diffusion model.

Finally, as we argued in the introduction, Boltzmann policies based on their own value function are attractive choices for pre-training of RL agents for later fine-tuning and multi-task settings. We hope to investigate this exciting potential in the future.

REFERENCES

Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B Tenenbaum, Tommi S Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh International Conference on Learning Representations*.

Tara Akhound-Sadegh, Jarrid Rector-Brooks, Avishek Joey Bose, Sarthak Mittal, Pablo Lemos, Cheng-Hao Liu, Marcin Sendera, Siamak Ravanbakhsh, Gauthier Gidel, Yoshua Bengio, Nikolay

- 486 Malkin, and Alexander Tong. Iterated denoising energy matching for sampling from boltzmann
487 densities. 2024.
488
- 489 Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow
490 network based generative models for non-iterative diverse candidate generation. *Advances in*
491 *Neural Information Processing Systems*, 34:27381–27394, 2021.
- 492 Yoshua Bengio, Salem Lahlou, Tristan Deleu, Edward J Hu, Mo Tiwari, and Emmanuel Bengio.
493 Gflownet foundations. *The Journal of Machine Learning Research*, 24(1):10006–10060, 2023.
494
- 495 Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy
496 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large
497 scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 498 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shu-
499 ran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint*
500 *arXiv:2303.04137*, 2023.
501
- 502 Tristan Deleu, Padideh Nouri, Nikolay Malkin, Doina Precup, and Yoshua Bengio. Discrete proba-
503 bilistic inference as control in multi-path environments. In *The 40th Conference on Uncertainty*
504 *in Artificial Intelligence*.
- 505 Zihan Ding and Chi Jin. Consistency models as a rich and efficient policy class for reinforcement
506 learning. In *The Twelfth International Conference on Learning Representations*, 2023.
507
- 508 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv*
509 *preprint arXiv:1605.08803*, 2016.
- 510 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with
511 deep energy-based policies. In *International conference on machine learning*, pp. 1352–1361.
512 PMLR, 2017.
- 513 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
514 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
515 *ence on machine learning*, pp. 1861–1870. PMLR, 2018a.
- 516 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
517 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and appli-
518 cations. *arXiv preprint arXiv:1812.05905*, 2018b.
519
- 520 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine.
521 Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint*
522 *arXiv:2304.10573*, 2023.
523
- 524 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
525 *neural information processing systems*, 33:6840–6851, 2020.
526
- 527 Vineet Jain and Siamak Ravanbakhsh. Learning to reach goals via diffusion. In *Forty-first Interna-*
528 *tional Conference on Machine Learning*, 2024.
- 529 Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for
530 flexible behavior synthesis. In *International Conference on Machine Learning*, pp. 9902–9915.
531 PMLR, 2022.
532
- 533 Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas.
534 Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*,
535 2021.
- 536 Bingyi Kang, Xiao Ma, Chao Du, Tianyu Pang, and Shuicheng Yan. Efficient diffusion policies for
537 offline reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
538
- 539 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
2014.

- 540 Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The*
541 *International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- 542 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-
543 learning. *arXiv preprint arXiv:2110.06169*, 2021.
- 544 Salem Lahlou, Tristan Deleu, Pablo Lemos, Dinghuai Zhang, Alexandra Volokhova, Alex
545 Hernández-García, Léna Néhale Ezzine, Yoshua Bengio, and Nikolay Malkin. A theory of con-
546 tinuous generative flow networks. In *International Conference on Machine Learning*, pp. 18269–
547 18300. PMLR, 2023.
- 548 Gene Li, Junbo Li, Anmol Kabra, Nati Srebro, Zhaoran Wang, and Zhuoran Yang. Exponential
549 family model-based reinforcement learning via score matching. *Advances in Neural Information*
550 *Processing Systems*, 35:28474–28487, 2022.
- 551 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Manfred Otto Heess, Tom Erez,
552 Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learn-
553 ing. *CoRR*, abs/1509.02971, 2015.
- 554 Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive energy
555 prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In *Inter-*
556 *national Conference on Machine Learning*, pp. 22825–22855. PMLR, 2023.
- 557 Kanika Madan, Jarrid Rector-Brooks, Maksym Korablyov, Emmanuel Bengio, Moksh Jain, An-
558 dreei Cristian Nica, Tom Bosc, Yoshua Bengio, and Nikolay Malkin. Learning gflownets from
559 partial episodes for improved convergence and stability. In *International Conference on Machine*
560 *Learning*, pp. 23467–23483. PMLR, 2023.
- 561 Nikolay Malkin, Moksh Jain, Emmanuel Bengio, Chen Sun, and Yoshua Bengio. Trajectory balance:
562 Improved credit assignment in gflownets. *Advances in Neural Information Processing Systems*,
563 35:5955–5967, 2022.
- 564 Michael Psenka, Alejandro Escontrela, Pieter Abbeel, and Yi Ma. Learning a diffusion model policy
565 from rewards via q-score matching. *arXiv preprint arXiv:2312.11752*, 2023.
- 566 Jarrid Rector-Brooks, Kanika Madan, Moksh Jain, Maksym Korablyov, Cheng-Hao Liu, Sarath
567 Chandar, Nikolay Malkin, and Yoshua Bengio. Thompson sampling for improved exploration in
568 gflownets. *arXiv preprint arXiv:2306.17693*, 2023.
- 569 Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned imitation
570 learning using score-based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 571 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon
572 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari,
573 go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 574 John Schulman, Sergey Levine, P. Abbeel, Michael I. Jordan, and Philipp Moritz. Trust region
575 policy optimization. *ArXiv*, abs/1502.05477, 2015.
- 576 Max W Shen, Emmanuel Bengio, Ehsan Hajiramezani, Andreas Loukas, Kyunghyun Cho, and
577 Tommaso Biancalani. Towards understanding and improving gflownet training. In *International*
578 *Conference on Machine Learning*, pp. 30956–30975. PMLR, 2023.
- 579 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
580 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
581 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 582 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
583 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-*
584 *ional Conference on Learning Representations (ICLR)*, 2021.
- 585 Niko Sünderhauf, Oliver Brock, Walter Scheirer, Raia Hadsell, Dieter Fox, Jürgen Leitner, Ben
586 Upcroft, Pieter Abbeel, Wolfram Burgard, Michael Milford, et al. The limits and potentials of
587 deep learning for robotics. *The International journal of robotics research*, 37(4-5):405–420, 2018.

594 Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
595
596 Daniil Tiapkin, Nikita Morozov, Alexey Naumov, and Dmitry P Vetrov. Generative flow networks
597 as entropy-regularized rl. In *International Conference on Artificial Intelligence and Statistics*, pp.
598 4213–4221. PMLR, 2024.

599 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy
600 class for offline reinforcement learning. In *The Eleventh International Conference on Learning
601 Representations*.
602

603 Philipp Wu, Alejandro Escontrela, Danijar Hafner, Pieter Abbeel, and Ken Goldberg. Daydreamer:
604 World models for physical robot learning. In *Conference on robot learning*, pp. 2226–2240.
605 PMLR, 2023.

606 Long Yang, Zhixiong Huang, Fenghao Lei, Yucun Zhong, Yiming Yang, Cong Fang, Shiting Wen,
607 Binbin Zhou, and Zhouchen Lin. Policy representation via diffusion probability model for rein-
608 forcement learning. *arXiv preprint arXiv:2305.13122*, 2023.

609 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse
610 reinforcement learning. In *Aaai*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A PROOFS

A.1 LEMMA 1 (POLICY IMPROVEMENT)

Lemma 1. *Let $\pi_{old} \in \Pi$ and π_{new} be the policy from Equation (6) with respect to $Q^{\pi_{old}}$. Then $Q^{\pi_{new}}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t)$ for all $\mathbf{s}_t, \mathbf{a}_t \in \mathcal{S}, \mathcal{A}$.*

Proof. Let $Q^{\pi_{old}}$ be the Q-function corresponding to π_{old} . Since π_{new} is defined using Equation (6) corresponding to $Q^{\pi_{old}}$, it must follow that,

$$D_{\text{KL}}(\pi_{new}(\cdot|s_t) || \exp(Q^{\pi_{old}}(s_t, \cdot) - \log Z^{\pi_{old}}(s_t))) \leq D_{\text{KL}}(\pi_{old}(\cdot|s_t) || \exp(Q^{\pi_{old}}(s_t, \cdot) - \log Z^{\pi_{old}}(s_t))),$$

where $D_{\text{KL}}(p||q)$ denotes the Kullback-Leibler divergence between two distributions p and q , and Z^π denotes the partition function $Z^\pi(s) = \int_a Q^\pi(s, a) da$.

Since the partition function depends only on the state, the inequality reduces to,

$$\begin{aligned} & \mathbb{E}_{a_t \sim \pi_{new}} [\log \pi_{new}(a_t|s_t) - Q^{\pi_{old}}(s_t, a_t) + \log Z^{\pi_{old}}(s_t)] \\ & \leq \mathbb{E}_{a_t \sim \pi_{old}} [\log \pi_{old}(a_t|s_t) - Q^{\pi_{old}}(s_t, a_t) + \log Z^{\pi_{old}}(s_t)] \\ \Rightarrow & \mathbb{E}_{a_t \sim \pi_{old}} [Q^{\pi_{old}}(s_t, a_t)] \leq \mathbb{E}_{a_t \sim \pi_{new}} [Q^{\pi_{old}}(s_t, a_t) - \log \pi_{new}(a_t|s_t) + \log \pi_{old}(a_t|s_t)]. \end{aligned} \quad (10)$$

Now, consider the non-negativity of the KL divergence between π_{new} and π_{old} ,

$$\begin{aligned} & D_{\text{KL}}(\pi_{new}(\cdot|s_t) || \pi_{old}(\cdot|s_t)) \geq 0 \\ \Rightarrow & \mathbb{E}_{a_t \sim \pi_{new}} [\log \pi_{new}(a_t|s_t) - \log \pi_{old}(a_t|s_t)] \geq 0 \end{aligned} \quad (11)$$

Combining Equation (12) and Equation (11), we get,

$$\mathbb{E}_{a_t \sim \pi_{old}} [Q^{\pi_{old}}(s_t, a_t)] \leq \mathbb{E}_{a_t \sim \pi_{new}} [Q^{\pi_{old}}(s_t, a_t)] \quad (12)$$

Now, consider the Bellman equation,

$$\begin{aligned} Q^{\pi_{old}}(s_t, a_t) &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_{t+1} \sim \pi_{old}} [Q^{\pi_{old}}(s_{t+1}, a_{t+1})]] \\ &\leq r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_{t+1} \sim \pi_{new}} [Q^{\pi_{old}}(s_{t+1}, a_{t+1})]] \quad (\text{from Equation (12)}) \\ &\leq r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_{t+1} \sim \pi_{new}} [r(s_{t+1}, a_{t+1}) + \gamma \mathbb{E}[\dots]]] \\ &\vdots \\ &\leq Q^{\pi_{new}}(s_t, a_t), \end{aligned} \quad (13)$$

where we repeatedly expand $Q^{\pi_{old}}$ in the RHS by applying the Bellman equation. The RHS converges to $Q^{\pi_{new}}$ following the fact that the Bellman operator is a contraction. \square

A.2 THEOREM 1 (POLICY ITERATION)

Theorem 1. *Alternating between policy evaluation and policy improvement repeatedly from any $\pi \in \Pi$ converges to a policy π^* such that $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ for all $\pi \in \Pi$ and $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{S}$.*

Proof. Let π_i be the policy at iteration i . By Lemma 1, Q^{π_i} is monotonically increasing. Since the reward is bounded, Q^π is also bounded for all $\pi \in \Pi$. This sequence converges to some policy π^* . At convergence, we must have for all $\pi \in \Pi$,

$$D_{\text{KL}}(\pi^*(\cdot|s_t) || \exp(Q^{\pi^*}(s_t, \cdot) - \log Z^{\pi^*}(s_t))) \leq D_{\text{KL}}(\pi(\cdot|s_t) || \exp(Q^{\pi^*}(s_t, \cdot) - \log Z^{\pi^*}(s_t))).$$

Using the same argument as in Lemma 1, we have $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ for all $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$, that is, the Q-value of any policy in Π is lower than that of π^* . Hence, π^* is optimal. \square

B IMPLEMENTATION DETAILS

The score function is parameterized as an MLP with two hidden layers of 256 units each with the ReLU activation function, except for the final layer. The MLP has skip connections as is typical for denoising score functions. The input to the policy comprises the state, noised action, the diffusion time step, and the temperature. The diffusion time step and the temperature are encoded using sinusoidal positional embeddings of 256 dimensions. The action is sampled following Equation (3) and the $\tanh(\cdot)$ function is applied to the sampled action followed by multiplication with the maximum value of the action space to ensure the value is within the correct range. The Q-network is also an MLP with two hidden layers of 256 units each with the ReLU activation function, except for the final layer. We use two Q-networks for the double Q-learning technique, and take the minimum of the two values.

The score function and the Q-network are trained for $250k$ environment steps with one mini-batch update per environment step. Optimization is performed using the Adam optimizer (Kingma, 2014) with a learning rate of 3×10^{-4} and a batch size of 256.

Table 1: Hyperparameters.

Parameter	Value
Number of hidden layers	2
Number of hidden units per layer	256
Sinusoidal embedding dimension	256
Activation function	ReLU
Optimizer	Adam
Learning rate	$3 \cdot 10^{-4}$
Number of samples per minibatch	256
Replay buffer size	250000
Discount factor	0.99
Gradient updates per step	1
Target smoothing co-efficient	0.005
Target update period	1
Seed training steps	10^4
σ_{\min}	0.00001
σ_{\max}	1
Number of Monte Carlo samples	1000
Number of integration steps	1000