

---

# Diffusion-based Curriculum Reinforcement Learning

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

Curriculum Reinforcement Learning (CRL) is an approach to facilitate the learning process of agents by structuring tasks in a sequence of increasing complexity. Despite its potential, many existing CRL methods struggle to efficiently guide agents toward desired outcomes, particularly in the absence of domain knowledge. This paper introduces DiCuRL (Diffusion Curriculum Reinforcement Learning), a novel method that leverages conditional diffusion models to generate curriculum goals. To estimate how close an agent is to achieving its goal, our method uniquely incorporates a  $Q$ -function and a trainable reward function based on Adversarial Intrinsic Motivation within the diffusion model. Furthermore, it promotes exploration through the inherent noising and denoising mechanism present in the diffusion models and is environment-agnostic. This combination allows for the generation of challenging yet achievable goals, enabling agents to learn effectively without relying on domain knowledge. We demonstrate the effectiveness of DiCuRL in three different maze environments simulated in MuJoCo, where it outperforms or matches nine state-of-the-art CRL algorithms from the literature.

## 1 Introduction

Reinforcement learning (RL) is a computational method that allows an agent to discover optimal actions through trial and error by receiving rewards and adapting its strategy to maximize cumulative rewards. Deep RL, which integrates deep neural networks (NNs) with RL, is an effective way to solve large-dimensional decision-making problems, such as learning to play video games [1, 2], chess [3], Go [4], and robot manipulation tasks [5, 6, 7, 8]. Without assuming knowledge about the environment, deep RL can tackle difficult uninformed search problems where the expected behaviors and rewards are often sparsely observed. To be effective though, deep RL typically requires a thorough exploration of the state space, which can be costly especially when the dimensionality of this space grows. Some methods, such as reward shaping [9], can mitigate the burden of exploration, but they require domain knowledge and prior task inspection, which makes these methods of limited applicability. Consequently, addressing the exploration efficiency in a domain-agnostic, algorithmic way has become a growing research interest. Different strategies, such as prioritizing replay sampling [8, 10, 11] or generating intermediate goals [12, 13, 14, 15, 16, 17, 18, 19], have been proposed to enhance efficiency. The latter approach, known as *curriculum learning*, focuses on designing a suitable curriculum to guide the agent gradually toward the desired goal. In generating curriculum goals, various approaches have been proposed. Some methods focus on interpolation between a source task distribution and a target task distribution [20, 21, 22, 17]. These methods rely on Kullback–Leibler (KL) divergence to measure the distance between distributions, which however has two main drawbacks: 1) it necessitates the distributions to be parameterized (e.g. as Gaussian distributions), and 2) it implicitly assumes Euclidean space, overlooking the manifold structure that characterizes some RL environments (for instance, in navigation problems where a direct path to the goal may be infeasible due to the presence of obstacles). Alternative approaches adopt optimal transport [13, 23]. Such approaches have mainly been applied to scenarios without significant exploration challenges, assuming the use of either

40 Wasserstein distance or geodesic interpolation between distribution sets. Curriculum generation based  
41 on uncertainty awareness has also been investigated [15, 24, 12]. Although [15] improves sample  
42 efficiency and generalization by selectively sampling training levels based on their estimated learning  
43 potential, it is limited by the need for uniformly sampling the remaining unseen levels. [24, 12]  
44 leverage meta-learning classifiers for uncertainty estimation, but struggle to distinguish uncertain  
45 areas as the goal space expands. Additional research tried to minimize the distance between the  
46 generated curriculum distribution and the desired outcome distribution, to propose intermediate task  
47 goals [19, 25]. However, these works use Euclidean distance, which as said can be problematic in  
48 some environments. To address this limitation, alternative methods [26, 27] incorporate graph-based  
49 planning, but they require the explicit specification of each obstacle’s boundary dimensions. Several  
50 studies [14, 28, 29] have proposed combining generative AI models, such as Generative Adversarial  
51 Networks (GANs) or diffusion models, with RL. In particular, [14] proposed a method for generating  
52 tasks with intermediate difficulty levels using GANs; however, it relies on an empirically chosen  
53 threshold. [28, 30, 29], instead, employ diffusion models within the context of offline RL, focusing  
54 on trajectory planning and policy optimization perspectives, respectively.

55 Despite these advancements, existing Curriculum Reinforcement Learning (CRL) approaches still  
56 struggle with generating suitable intermediate goals, particularly in complex environments with  
57 significant exploration challenges.

58 To overcome this challenge, in this paper, we propose DiCURL (Diffusion Curriculum Reinforcement  
59 Learning). Our method leverages conditional diffusion models to dynamically generate curriculum  
60 goals, guiding agents towards desired goals while simultaneously considering the  $Q$ -function and  
61 a trainable reward function based on Adversarial Intrinsic Motivation (AIM) [31]. Unlike previous  
62 offline RL approaches [28, 30, 29] that train and use diffusion models for planning or policy generation  
63 relying on pre-existing data, DiCURL facilitates online learning and adaptation by structuring  
64 achievable curriculum goals using a diffusion model within a CRL setting. This enables agents to  
65 learn effectively without requiring domain-specific knowledge. Our approach achieves this by the  
66 following key steps: ① The diffusion model captures the distribution of visited states and facilitates  
67 exploration through its inherent noising and denoising mechanism. ② As the  $Q$ -function predicts the  
68 cumulative reward starting from a state and a given goal while following a policy, we can determine  
69 feasible goals by maximizing the  $Q$ -function, ensuring that the generated goals are challenging yet  
70 achievable for the agent. ③ The AIM reward function estimates the agent’s proximity to the desired  
71 goal and allows us to progressively shift the curriculum towards the desired goal.

72 We compare our proposed approach with nine state-of-the-art CRL baselines in three different maze  
73 environments simulated in MuJoCo [32]. Our results show that DiCURL surpasses or performs on  
74 par with the state-of-the-art CRL algorithms. To summarize, our main contributions are as follows:

- 75 • We present DiCURL, a novel approach to CRL that leverages a diffusion model to accurately  
76 capture the visited state distribution and facilitate exploration through noising and denoising.  
77 • Our novel integration of the  $Q$ -function and the AIM reward function into the diffusion model  
78 allows the generation of curriculum goals that are optimally challenging for the agent, while  
79 progressively guiding it toward the desired goal during training.

## 80 2 Related Work

81 CRL [33] algorithms generally adjust the sequence of learning experiences to improve the agent’s  
82 performance or accelerate training. These algorithms focus on formulating intermediate goals that  
83 progressively guide the agent toward the desired goal.

84 Hindsight Goal Generation (HGG) [19] addresses the inefficiency issue inherent in Hindsight Experi-  
85 ence Replay (HER) [34] by generating hindsight goals through maximizing a value function and  
86 minimizing the Wasserstein distance between the achieved goal and the desired goal distribution. The  
87 Wasserstein distance is approximated using the Euclidean distance, which as said is not suitable in  
88 environments such as mazes where walls obstruct the direct path between the initial and desired goal.  
89 On the other hand, in obstacle-free environments, this approach enables HGG to select short-term  
90 achievable goals that effectively guide the agent toward the desired goal.

91 CURROT [23] and GRADIENT [13] both employ optimal transport for the generation of inter-  
92 mediate goals. CURROT formulates CRL as a constrained optimization problem and uses the

93 Wasserstein distance to measure the distance between distributions. Conversely, GRADIENT intro-  
94 duces task-dependent contextual distance metrics and can manage non-parametric distributions in  
95 both continuous and discrete context settings; moreover, it directly interprets the interpolation as the  
96 geodesic from the source to the target distribution.

97 GOAL-GAN [14] generates intermediate goals using a Generative Adversarial Network (GAN) [35],  
98 without considering the target distribution. A goal generator is used to propose goal regions, and  
99 a goal discriminator is trained to evaluate if a goal is at the right level of difficulty for the current  
100 policy. The specification of goal regions is done using an indicator reward function, and policies are  
101 conditioned on the goal as well as the state, similarly to a universal value function approximator [36].

102 PLR [15] uses selective sampling to prioritize instances with higher estimated learning potential  
103 for future revisits during training. Learning potential is estimated using TD-Errors, resulting in the  
104 creation of a more challenging curriculum. VSD [16] estimates the epistemic uncertainty of the value  
105 function and selects goals based on this uncertainty measure. The value function confidently assigns  
106 high values to easily achievable goals and low values to overly challenging ones. ACL [18] maximizes  
107 learning progress by considering two main measures: the rate of improvement in prediction accuracy  
108 and the rate of increase in network complexity. These indicators, while explaining the current rate  
109 of improvement of the learner, act as a reward signal and are further provided to a non-stationary  
110 multi-armed bandit algorithm that determines the curriculum. The ALP-GMM [17] fits Gaussian  
111 Mixture Models (GMM) using an Absolute Learning Progress (ALP) score, which is defined as  
112 the absolute difference in rewards between the current episode and a previous episode. The teacher  
113 algorithm generates curriculum goals by sampling environments to maximize the student’s ALP,  
114 which is modeled by the GMM.

115 OUTPACE [12] employs a trainable intrinsic reward mechanism, known as Adversarial Intrinsic  
116 Motivation (AIM) [31] (the same used in our method) which is designed to minimize the Wasserstein  
117 distance between the state visitation distribution and the goal distribution. This function increases  
118 along the optimal goal-reaching trajectory. For curriculum goal generation, OUTPACE uses Con-  
119 ditional Normalized Maximum Likelihood (CNML), to classify state success labels based on their  
120 association with visited states, out-of-distribution samples, or the desired goal distribution. Finally, it  
121 prioritizes uncertain and temporally distant goals using meta-learning-based uncertainty quantification  
122 [37] and Wasserstein-distance-based temporal distance approximation.

### 123 3 Background

124 We now introduce the background concepts on multi-goal RL, Soft Actor-Critic (SAC), Wasserstein  
125 distance, Adversarial Intrinsic Motivation (AIM), and diffusion models.

#### 126 3.1 Multi-Goal Reinforcement Learning

127 In the context of multi-goal RL, we can formulate the RL problem as a goal-oriented Markov decision  
128 process (MDP) characterized by continuous state and action spaces. This MDP is defined by the  
129 tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{G}, \mathcal{T}, p, \gamma_r \rangle$ . Here,  $\mathcal{S}$  represents the state space,  $\mathcal{A}$  denotes the action space, and  $\mathcal{G}$  is the  
130 goal space. The transition dynamics,  $\mathcal{T}(s'|s, a)$ , describes the probability of transitioning to the next  
131 state  $s'$  given the current state  $s$  and action  $a$ . The joint probability distribution over the initial states  
132 and the desired goal distribution is represented by  $p(s_0, g)$ , and  $\gamma_r \in [0, 1]$  is a discount factor.

133 We utilize the AIM reward function  $r$ , as outlined in Section 3.3. The objective is to identify  
134 the optimal policy  $\pi$  that maximizes the expected cumulative reward approximated by the value  
135 function  $Q^\pi(s, g, a)$ . This function can be expanded to the Universal Value Function (UVF), a  
136 goal-conditioned value function incorporating the goal into value estimation. The UVF, defined as  
137  $V^\pi(s, g)$ , where  $s$  is the current state,  $g$  is the goal, and  $\pi$  is the policy, estimates the expected return  
138 from state  $s$  under policy  $\pi$ , given goal  $g$ . The UVF has been shown to successfully generalize to  
139 unseen goals, which makes it a viable approach for multi-goal RL [36].

#### 140 3.2 Soft Actor-Critic

141 Soft Actor-Critic (SAC) [38] is an off-policy RL algorithm that learns a policy maximizing the sum  
142 of the expected discounted cumulative reward and the entropy of the policy, namely:  $J(\pi) =$

143  $\sum_{t=0}^T \mathbb{E}_{s_t \sim \mathcal{B}} r + \alpha_e \mathcal{H}(\pi(\cdot | s_t, g))$ . The policy  $\pi_\psi(a | s) : \mathcal{S} \mapsto \mathcal{P}(\mathcal{A})$  defines a map from  
 144 state to distributions over actions, parameterized by  $\psi$ . A state-action value function is defined  
 145 as  $Q_\phi(s, g, a) : \mathcal{S} \times \mathcal{G} \times \mathcal{A} \rightarrow \mathbb{R}$ , parameterized by  $\phi$ . The policy parameters can be learned by  
 146 directly minimizing the expected KL divergence and the  $Q$ -value is trained by minimizing the soft  
 147 Bellman residual, which are defined as the following loss functions:

$$\mathcal{L}_\pi = \mathbb{E}_{(s,g) \sim \mathcal{B}} [\mathbb{E}_{a \sim \pi_\psi} [\alpha_e \log (\pi_\psi(a | s, g) - Q_\phi(s, g, a))] ] \quad (1)$$

$$\mathcal{L}_Q = \mathbb{E}_{(s,a,r,s',g) \sim \mathcal{B}} [Q_\phi(s, g, a) - (r + \gamma_r \mathbb{E}_{s' \sim \mathcal{B}} [V_\phi(s', g)])^2] \quad (2)$$

148 where  $V_\phi(s, g) = \mathbb{E}_{a \sim \pi} [Q_\phi(s, g, a) - \log \pi_\psi(a | s, g)]$ ,  $\mathcal{B}$  is the replay buffer and  $\alpha_e$  is the temper-  
 149 ature parameter that controls the stochasticity of the optimal policy.

### 150 3.3 Wasserstein Distance and Adversarial Intrinsic Motivation Reward Function

151 The Wasserstein distance offers a way to quantify the amount of work required to transport one  
 152 distribution to another distribution. The Wasserstein-p distance  $W_p$  between two distributions  $\mu$  and  
 153  $\nu$  on a metric space  $(\mathcal{X}, d)$ , where  $\mathcal{X}$  is a set and  $d$  denotes a metric on  $\mathcal{X}$ , is defined as follows [39]:

$$W_p(\mu, \nu) := \inf_{\gamma \in \Pi(\mu, \nu)} \left( \int_{\mathcal{X}} d(x, y)^p d\gamma(x, y) \right)^{1/p} = \inf_{\gamma \in \Pi(\mu, \nu)} \mathbb{E}_{(X,Y) \sim \gamma} [d(X, Y)^p]^{1/p} \quad (3)$$

154 where  $\Pi(\mu, \nu)$  denotes the set of all possible joint distributions  $\gamma(x, y)$  whose marginals are  $\mu$  and  
 155  $\nu$ . Intuitively,  $\gamma(x, y)$  tells what is the least amount of work, as measured by  $d$ , that needs to be  
 156 done to convert the distribution  $\mu$  into the distribution  $\nu$  [40]. A timestep quasi-metric  $d^\pi(s, g)$  can  
 157 be used as a distance metric. It estimates the work needed to transport one distribution to another,  
 158 representing the number of transition steps required to reach the goal state  $g \in \mathcal{G}$  for the first time  
 159 when following the policy  $\pi$ .

160 As proposed by Durugkar et al. [31], the AIM reward function can be learned by minimizing the  
 161 Wasserstein distance between the state visitation distribution  $\rho_\pi$  and the desired goal distribution  
 162  $\mathcal{G}$ . Through the minimization of the Wasserstein distance  $W_1(\rho_\pi, \mathcal{G})$  (for  $p = 1$ ,  $W_1$  is known as  
 163 the Kantorovich–Rubinstein distance), a reward function can be formulated to estimate the work  
 164 required to transport the state visitation distribution  $\rho_\pi$  to the desired goal distribution  $\mathcal{G}$  and is  
 165 expressed as follows:  $W_1(\rho_\pi, \mathcal{G}) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{g \sim \mathcal{G}} [f(g)] - \mathbb{E}_{s \sim \rho_\pi} [f(s)]]$ . If the state visitation  
 166 distribution  $\rho_\pi(s)$  comprises states that optimally progress towards the goal  $g$ , the potential function  
 167  $f(s)$  increases along the trajectory, reaching its maximum value at  $f(g)$ . The reward function, which  
 168 can be approximated using a neural network, denoted as  $r_\varphi^\pi$ , increases as the states approach the  
 169 desired goal  $g \in \mathcal{G}$ .  $r_\varphi^\pi$  can be trained using the data collected by the policy  $\pi$ . Leveraging the  
 170 estimation of the Wasserstein distance  $W_1(\rho_\pi, \mathcal{G})$ , the loss function for training the parameterized  
 171 reward function  $r_\varphi^\pi$  is defined as follows:

$$\mathcal{L}_\varphi = \mathbb{E}_{(s,g) \sim \mathcal{B}} [f_\varphi^\pi(s) - f_\varphi^\pi(g)] + \lambda \cdot \mathbb{E}_{(s,s',g) \sim \mathcal{B}} [\max(|f_\varphi^\pi(s) - f_\varphi^\pi(s')| - 1, 0)]^2 \quad (4)$$

172 where the second component of the sum is a penalty term and the coefficient  $\lambda$  is necessary to ensure  
 173 smoothness [31]. The reward can be calculated as  $r_\varphi(s, g) = f_\varphi^\pi(s) - f_\varphi^\pi(g)$ , which is the negative  
 174 of the Wasserstein distance  $-W_1(\rho_\pi, \mathcal{G})$ .

### 175 3.4 Diffusion Models

176 Diffusion models [41] express a probability distribution  $p(x_0)$  through latent variables in the form  
 177  $p_\theta(\mathbf{x}_0) := \int p_\theta(\mathbf{x}_{0:N}) d\mathbf{x}_{1:N}$ , where  $\mathbf{x}_1, \dots, \mathbf{x}_N$  are latent variables of the same dimensionality  
 178 as the data  $\mathbf{x}_0 \sim p(\mathbf{x}_0)$ . They are characterized by a forward and a reverse diffusion process.  
 179 The forward diffusion process approximates the posterior  $q(\mathbf{x}_{1:N} | \mathbf{x}_0)$  using a Markov chain that  
 180 perturbs the input data by gradually adding Gaussian noise  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$  in  $N$  steps with a predefined  
 181 variance schedule  $\beta_1, \dots, \beta_N$ . This process is defined as:

$$q(\mathbf{x}_{1:N} | \mathbf{x}_0) := \prod_{k=1}^N q(\mathbf{x}_k | \mathbf{x}_{k-1}) \quad (5)$$

$$q(\mathbf{x}_k | \mathbf{x}_{k-1}) := \mathcal{N}(\mathbf{x}_k; \sqrt{1 - \beta_k} \mathbf{x}_{k-1}, \beta_k \mathbf{I})$$

182 The reverse diffusion process aims to recover the original input data from the noisy (diffused) data.  
 183 It learns to progressively reverse the diffusion process step by step and approximates the joint  
 184 distribution  $p_\theta(\mathbf{x}_{0:N})$ . This process is defined as:

$$p_\theta(\mathbf{x}_{0:N}) := p(\mathbf{x}_N) \prod_{k=1}^N p_\theta(\mathbf{x}_{k-1} | \mathbf{x}_k) \quad (6)$$

$$p_\theta(\mathbf{x}_{k-1} | \mathbf{x}_k) := \mathcal{N}(\mathbf{x}_{k-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_k, k), \boldsymbol{\Sigma}_\theta(\mathbf{x}_k, k))$$

185 where  $p(\mathbf{x}_N) = \mathcal{N}(\mathbf{x}_N; \mathbf{0}, \mathbf{I})$ . The optimization of the reverse diffusion process is achieved by  
 186 maximizing the evidence lower bound (ELBO)  $\mathbb{E}_q \left[ \ln \frac{p_\theta(\mathbf{x}_{0:N})}{q(\mathbf{x}_{1:N} | \mathbf{x}_0)} \right]$  [42]. Once trained, sampling data  
 187 from Gaussian noise  $\mathbf{x}_N \sim p(\mathbf{x}_N)$  and running through the reverse diffusion process from  $k = N$  to  
 188  $k = 0$  yields an approximation of the original data distribution.

## 189 4 Methodology

190 As discussed earlier, in multi-goal RL, the desired goal is sampled from a desired goal distribution  
 191 in each episode, and the agent aims to achieve multiple goals. By integrating a curriculum design  
 192 into multi-goal RL, we can reformulate a task in such a way that it starts with easier goals and  
 193 progressively increases in difficulty. Our curriculum diffusion-model-based goal-generation method  
 194 works as follows. Given the presence of two different types of timesteps for the diffusion process  
 195 and the RL task, we denote the diffusion timesteps using subscripts  $k \in \{1, \dots, N\}$  and the RL  
 196 trajectory timesteps using subscripts  $t \in \{1, \dots, T\}$ . Our curriculum goal generation takes the state  
 197  $s$  as an input. It then outputs a curriculum goal set  $\mathcal{G}_c$ , which is obtained through the reverse diffusion  
 198 process of a conditional diffusion model as follows:

$$\mathcal{G}_c = p_\theta(\mathbf{g}_{0:N} | s) = \mathcal{N}(\mathbf{g}_N; \mathbf{0}, \mathbf{I}) \prod_{k=1}^N p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k, s) \quad (7)$$

199 where  $\mathbf{g}_0$  is the end sample of the reverse diffusion process used as a curriculum goal. Commonly,  
 200  $p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k, s)$  is a conditional distribution parametrized by  $\theta$  and is chosen to model a multivariate  
 201 Gaussian distribution  $\mathcal{N}(\mathbf{g}_{k-1}; \boldsymbol{\mu}_\theta(\mathbf{g}_k, s, k), \boldsymbol{\Sigma}_\theta(\mathbf{g}_k, s, k))$ . Following [41], rather than learning  
 202 the variances of the forward diffusion process, we assign a fixed covariance matrix  $\boldsymbol{\Sigma}_\theta(\mathbf{g}_i, s, i) = \beta_i \mathbf{I}$   
 203 and a mean defined as:

$$\boldsymbol{\mu}_\theta(\mathbf{g}_k, s, k) = \frac{1}{\sqrt{\alpha_k}} \left( \mathbf{g}_k - \frac{\beta_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(\mathbf{g}_k, s, k) \right). \quad (8)$$

204 Initially, we sample a Gaussian noise  $\mathbf{g}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . We then apply the reverse diffusion process  
 205 parameterized by  $\theta$  starting from the last step  $N$  and proceeding backward to step 1:

$$\mathbf{g}_{k-1} | \mathbf{g}_k = \frac{\mathbf{g}_k}{\sqrt{\alpha_k}} - \frac{\beta_k}{\sqrt{\alpha_k(1 - \bar{\alpha}_k)}} \epsilon_\theta(\mathbf{g}_k, s, k) + \sqrt{\beta_k} \epsilon \quad (9)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \text{ for } k = N, \dots, 1.$$

206 For the case  $k = 1$  (last term of the reverse diffusion process), we set  $\epsilon$  to  $\mathbf{0}$ , to enhance the sampling  
 207 quality by ignoring  $\sqrt{\beta_k}$  and edge effects. In fact, as empirically demonstrated by [41], training  
 208 the diffusion model yields better results when utilizing a simplified loss function that excludes  
 209 the weighting term. Thus, we adopt the following simplified loss function to train the conditional  
 210  $\epsilon_\theta$ -model, where  $\epsilon_\theta$  is a function approximator designed to predict  $\epsilon^1$ .

$$\mathcal{L}_d(\theta) = \mathbb{E}_{k \sim \mathcal{U}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), (s, \mathbf{g}_0) \sim \mathcal{B}} \left[ \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_k} \mathbf{g}_0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, s, k) \right\|^2 \right], \quad (10)$$

211 where  $\mathcal{U}$  is a uniform distribution over the discrete set  $\{1, \dots, N\}$  and  $\mathcal{B}$  denotes the replay buffer  
 212 collected by policy  $\pi$ .

---

<sup>1</sup>Details about deriving the conditional diffusion model equation are provided in Supp. Mat. A.

213 To sample feasible curriculum goals for the agent, we integrate the  $Q$ -function and the AIM reward  
 214 function into the loss. The total loss function for training the diffusion model is defined as follows:

$$\mathcal{L} = \xi_d \cdot \mathcal{L}_d(\theta) - \mathbb{E}_{s,a \sim \mathcal{B}, \mathbf{g}_0 \sim \mathcal{G}_c, g_d \sim \mathcal{G}} [\xi_q \cdot Q_\phi(s, \mathbf{g}_0, \pi(s, \mathbf{g}_0)) + \xi_r \cdot r_\varphi(\mathbf{g}_0, g_d)] \quad (11)$$

215 The rationale is that, while minimizing the loss  $\mathcal{L}_d$  allows us to accurately capture the state distribution,  
 216 we also aim to simultaneously maximize the expected value of  $Q$  and the AIM reward function  $r$ ,  
 217 using the weights  $\xi_d$ ,  $\xi_q$ , and  $\xi_r$  to adjust the relative importance of the three components of the loss  
 218 function. More in detail, the  $Q$ -function predicts the cumulative reward starting from a state and  
 219 following the policy, while the AIM reward function estimates how close an agent is to achieving  
 220 its goal. By maximizing  $Q$  and the AIM reward, we can generate curriculum goals that are neither  
 221 overly simplistic nor excessively challenging, progressing towards the desired goal.

222 In Eq. (11),  $\mathbf{g}_0$  is obtained by sampling experiences from the replay buffer using Eq. (9) through  
 223 a reverse diffusion process parameterized by  $\theta$ . Therefore, taking the gradient of  $Q$  and the AIM  
 224 reward  $r$  involves back-propagating through the entire diffusion process. After we obtain the set of  
 225 curriculum goals  $\mathcal{G}_c$  from the diffusion model, we use a bipartite graph  $G(\{V_x, V_y\}, E)$  with edge  
 226 costs  $w$ , composed of vertices  $V_x$  (i.e., the curriculum goal candidates derived from the diffusion  
 227 model) and vertices  $V_y$  (i.e., the desired goals), where  $E$  represents the edge weights, and select the  
 228 optimal curriculum goal  $g_c^*$ <sup>2</sup>. We employ the Minimum Cost Maximum Flow algorithm in order to  
 229 solve the bipartite matching problem and identify the edges with the minimal cost  $w$ :

$$\max_{\hat{\mathcal{G}}^c: |\hat{\mathcal{G}}^c|=K} \sum_{\substack{g_t^0, \dots, T \in \hat{\mathcal{G}}^d, g_+^i \in \hat{\mathcal{G}}^+}} w \quad \text{where: } w = \sqrt{\frac{g_i - \bar{g}}{N}}. \quad (12)$$

230 The specifics of the generation of intermediate goals through diffusion models are explained in  
 231 Algorithm 1, while the overall algorithm is outlined in Algorithm 2.

---

### Algorithm 1 Diffusion Curriculum Goal Generator

---

```

1: Input: state  $s$ , no. of reverse diffusion timestep  $N$ , training step  $M$ 
2: Obtain states  $s$  from the minibatch  $b$ 
3: for  $i = 1, \dots, M$  do ▷ Training iterations
4:    $\mathbf{g}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   for  $k = N, \dots, 1$  do ▷ Reverse diffusion process
6:      $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
7:      $g_{k-1} = \frac{1}{\sqrt{\alpha_k}} \left( g_k - \frac{\beta_k}{\sqrt{1-\alpha_k}} \epsilon_\theta(g_k, k, s) \right) + \sqrt{\beta_k} \epsilon$  ▷ Using Eq. (9)
8:      $k \sim \text{Uniform}(\{1, \dots, N\})$ 
9:      $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
10:    Calculate diffusion  $\mathcal{L}_d(\theta), Q(s, \mathbf{g}_0, \pi(s)), r(\mathbf{g}_0, g_d)$  ▷ Calculate with the generated goal  $\mathbf{g}_0$ 
11:    Calculate total loss  $\mathcal{L} = \xi_d \mathcal{L}_d(\theta) - \xi_q Q(s, \mathbf{g}_0, \pi(s)) - \xi_r r(\mathbf{g}_0, g_d)$  ▷ Eq 11
12:     $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}$  ▷ Take gradient descent step and update the diffusion weights
return  $\mathbf{g}_0$ 

```

---

232 In Algorithm 1, during the training iterations, we sample Gaussian noise  $\mathbf{g}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  to denoise  
 233 the data in the reverse diffusion process. Between lines 5 and 7, we iterate from time step  $N$  to 1,  
 234 where we perform the reverse diffusion process using Eq. (9) and subtract the predicted noise  $\epsilon_\theta$   
 235 from  $\mathbf{g}_N$  to denoise the noisy goal iteratively. We then uniformly sample a timestep  $k$  from the range  
 236 between 1 and  $N$  (line 8) and sample a Gaussian noise  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  (line 9) in order to calculate the  
 237 diffusion loss defined in Eq. (10). In line 10, we calculate the diffusion loss, the  $Q$ -value, and the  
 238 AIM reward function using the generated goal  $\mathbf{g}_0$  from the reverse diffusion process. Then in line  
 239 11, we calculate the total loss defined in Eq. (11) and update the diffusion model parameters  $\theta$  using  
 240 gradient descent. Finally, we return the generated goal  $\mathbf{g}_0$ .

<sup>2</sup>In principle, instead of sampling experiences from the replay buffer and providing them to the diffusion model to generate a curriculum goal set, only the state of the last timestep,  $s_T$ , could be provided to the diffusion model to obtain a curriculum goal. The difference is that giving the sampled experiences to the diffusion model generates many curriculum points, while using only the last timestep  $s_T$  generates only a single curriculum point, eliminating the need to apply any selection strategy. This is investigated in Supp. Mat. B.

---

**Algorithm 2** RL Training

---

```
1: Input: no. of episodes  $E$ , timesteps  $T$ 
2: Select an off-policy algorithm  $\mathbb{A}$  ▷ In our case,  $\mathbb{A}$  is SAC
3: Initialize replay buffer  $\mathcal{B} \leftarrow \emptyset$ ,  $g_c \leftarrow \{g_d\}$  and networks  $Q_\phi, \pi_\psi, r_\varphi$ 
4: for episode = 0 ...  $E$  do
5:   Sample initial state  $s_0$ 
6:   for  $t = 0 \dots T$  do
7:      $a_t = \pi(s_t, g_c)$ 
8:     Execute  $a_t$ , obtain next state  $s_{t+1}$ 
9:     Store transition  $(s_t, a_t, r_t, s_{t+1}, g_c)$  in  $\mathcal{B}$ 
10:    Sample a minibatch  $b$  from replay buffer  $\mathcal{B}$ 
11:     $\mathcal{G}_c \leftarrow \text{DiffusionCurriculumGenerator}(b)$ 
12:    Find  $g_c$  that maximizes  $w$  in Eq. (12)
13:    Update  $Q$  and  $\pi$  with  $b$  to minimize  $\mathcal{L}_Q$  and  $\mathcal{L}_\pi$  in Eq. (1) and in Eq. (2)
14:    Update the AIM reward function  $r_\varphi$ 
15:    $success \leftarrow 0$  ▷ Success rate
16:   Sample a desired goal  $g_d \sim \mathcal{G}$ 
17:   for  $i = 1 \dots n_{testrollout}$  do
18:      $a_t = \pi(s_t, g_d)$ 
19:     Execute  $a_t$ , obtain next state  $s_{t+1}$  and reward  $r_t$ 
20:     if  $|\phi(s_{t+1}) - g_d| \leq \kappa$  then
21:        $success = success + 1/n_{testrollout}$ 
```

---

241 In Algorithm 2 we begin by defining an off-policy algorithm denoted as  $\mathbb{A}$ . While any off-policy  
242 algorithm could be employed, we choose Soft Actor-Critic (SAC) to align with the baseline algorithms  
243 tested in our experiments. In line 5, we sample the initial state and provide the curriculum goal  $g_c$  to  
244 the policy, along with the current state. The policy generates an action (line 7), and executes it in the  
245 environment (line 8). Subsequently, the next state and reward are obtained (line 9). Then, we provide  
246 the minibatch  $b$  to the curriculum goal generator in line 11 to generate a curriculum goal set  $\mathcal{G}_c$ . Then  
247 we find the optimal curriculum goal  $g_c$  using bipartite graph optimization (line 12). Furthermore, the  
248 loss functions defined in Eq. (1) and Eq. (2) are calculated, and the networks approximating  $\pi$  and  $Q$ ,  
249 as well as the AIM reward function  $r$  defined in Eq. (4), are updated. Between line 15 and 21, we run  
250  $n$  test rollouts and extract the achieved state of the agent using  $\phi$  and calculate the success rate in  
251 reaching the desired goal within a threshold value  $\kappa$ .

252 **5 Experiments**

253 To evaluate our proposed method, we conducted experiments across three maze environments  
254 simulated in MuJoCo<sup>3</sup>: *PointUMaze*, *PointNMaze*, and *PointSpiralMaze*. In these environments, a  
255 goal is interpreted as the  $(x, y)$  position of the agent achieved in an episode. These environments  
256 have been specifically chosen due to their structural characteristics, which are ideal for testing  
257 environment-agnostic curriculum generation strategies. Moreover, they present a variety of complex  
258 and diverse navigation challenges, requiring an agent to learn effective exploration and exploitation.

259 We compared our approach DiCURL against nine state-of-the-art CRL algorithms, namely ACL [18],  
260 GOAL-GAN [14], HGG [19], ALP-GMM [17], VDS [16], PLR [15], CURROT [23], GRADIENT [13],  
261 and OUTPACE [12], each one run with 5 different random seeds. The primary conceptual  
262 differences between ours and baseline algorithms are summarized in Table 1. Details on the baseline  
263 algorithms, parametrization, training setup, and maze environments are given in Supp. Mat. C. Our  
264 codebase is available at: <https://anonymous.4open.science/r/diffusioncurriculum/>.

265 The results, shown in Fig. 1 and detailed in Table 2<sup>4</sup>, demonstrate the effectiveness of the proposed  
266 DiCURL method. Notably, DiCURL outperforms or matches all the baseline methods. For the  
267 PointNMaze and PointSpiralMaze, the success rate of all methods, except for ours and OUTPACE  
268 (and HGG for PointNMaze), is close to zero (detailed results are reported in Supp. Mat. C). ACL,

<sup>3</sup>Details available at [https://robotics.farama.org/envs/maze/point\\_maze/](https://robotics.farama.org/envs/maze/point_maze/).

<sup>4</sup>We only show methods that reach success rate 1.0 within the allotted timesteps on at least one environment.

Table 1: Comparison of DiCURL with previous CRL methods from the literature (sorted by year).

Algorithm	Curriculum method	Target dist. curriculum	Geometry-agnostic	Off-policy	External reward	Venue, year
ACL [18]	LSTM	X	✓	X	X	PMLR, 2017
GoalGAN [14]	GAN	X	✓	X	X	PMLR, 2018
HGG [19]	$Q, \mathcal{B}, W_2$	$\mathcal{G}^+$	X	✓	X	NeurIPS, 2019
ALP-GMM [17]	GMM	X	✓	✓	X	PMLR, 2020
VDS [16]	$Q, \mathcal{B}$	X	✓	✓	X	NeurIPS, 2020
PLR [15]	TD-Error	X	✓	X	X	PMLR, 2021
CURROT [23]	$W_2$	$\mathcal{G}^+, \mathcal{U}$	X	✓	X	PMLR, 2022
GRADIENT [13]	$W_2$	$\mathcal{G}^+$	X	✓	✓	NeurIPS, 2022
OUTPACE [12]	CNML	$\mathcal{G}^+$	✓	✓	✓	ICLR, 2023
DiCURL (Ours)	Diffusion	$\mathcal{G}^+$	✓	✓	✓	

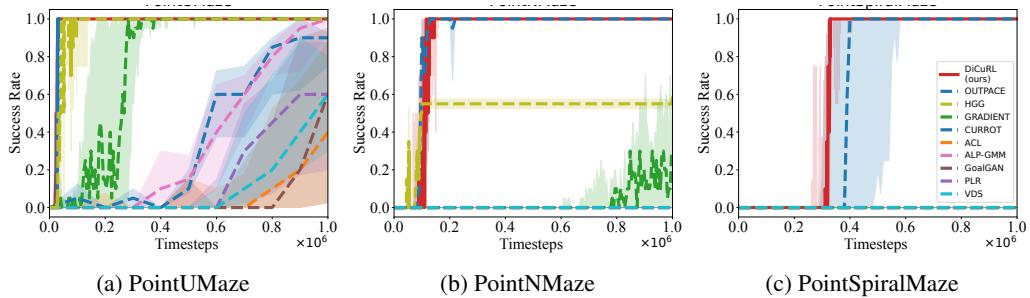


Figure 1: Test success rate for the algorithms under comparison on the three maze tasks.

Table 2: No. of timesteps (rounded) to reach success rate 1.0 (average  $\pm$  std. dev. across 5 runs). NaN indicates that a success rate of 1.0 is not reached within the maximum budget of 1e6 timesteps.

Algorithm	PointUMaze	PointNMaze	PointSpiralMaze
DiCURL (Ours)	$28333 \pm 3036$	$108428 \pm 34718$	$305833 \pm 43225$
OUTPACE	$29800 \pm 4166$	$113333 \pm 24267$	$396875 \pm 111451$
HGG	$48750 \pm 24314$	NaN	NaN
GRADIENT	$263431 \pm 114795$	NaN	NaN

269 GoalGAN, ALP-GMM, VDS, and PLR, which lack awareness of the target distribution, underperform compared to target-aware methods such as our proposed approach, OUTPACE, GRADIENT,  
270 CURROT, and HGG. HGG encounters difficulties because of infeasible curriculum proposals. This  
271 is a result of its reliance on the Euclidean distance metric.

272 Both CURROT and GRADIENT generate curriculum goals using the optimal transport method and  
273 the Wasserstein distance metric, relying on the geometry of the environment. This dependence might  
274 be the reason for their inconsistent and poor performance across the different environments.

275 OUTPACE, utilizes instead the Wasserstein distance and an uncertainty classifier for uncertainty-  
276 aware curriculum learning, exhibiting similar performance to our approach but with slower conver-  
277 gence and higher variance in success rate. This is likely due to the fact that OUTPACE's curriculum  
278 heavily relies on the visited state distributions, which necessitate an initial exploration by the agent.  
279 While also our approach depends on visited states, incorporating the Qand AIM reward functions into  
280 curriculum generation facilitates exploration beyond them, potentially explaining the performance  
281 differences between DiCURL and OUTPACE.

282 Lastly, we recall that our approach generates curriculum goals based on the reverse diffusion process,  
283 which allows the generated curriculum goals to gradually shift from the initial state distribution to the  
284 desired goal distribution. Fig. 2 shows an example of a curriculum set generated by DiCURL, for  
285 the case of the PointUMaze environment, at each iteration of the reverse diffusion process. Fig. 3,  
286 instead, shows the differences between the curriculum goals generated by DiCURL, GRADIENT,  
287 and HGG in the PointSpiralMaze environment: it can be seen how DiCURL manages to generate  
288 curriculum goals that explore the whole environment more effectively than the baseline algorithms.

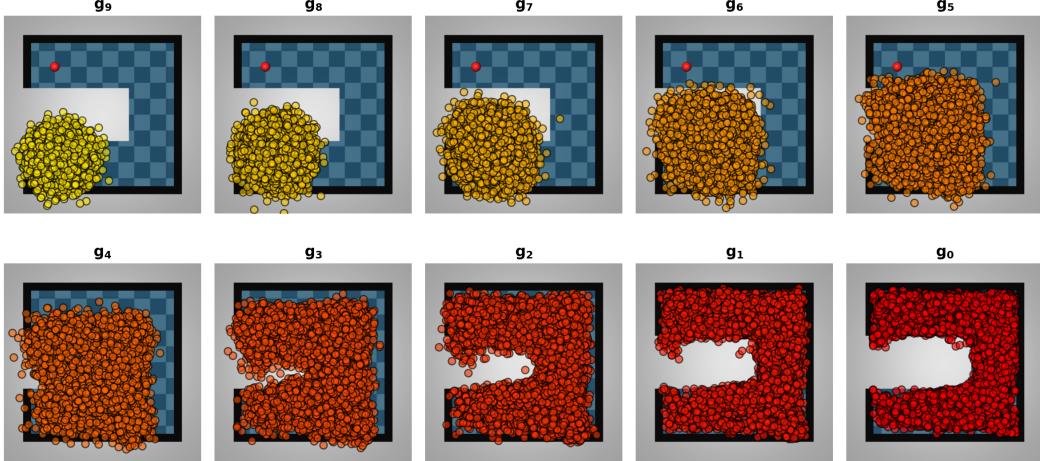


Figure 2: The curriculum goal set  $\mathcal{G}_c$  generated by DiCURL during the reverse diffusion process (lines 5-7 in Algorithm 1) for the PointUMaze environment. The color indicates the goals generated at a specific iteration step during the reverse diffusion process of the diffusion model. Then these goals are selected in Eq. (12) based on the given cost function.

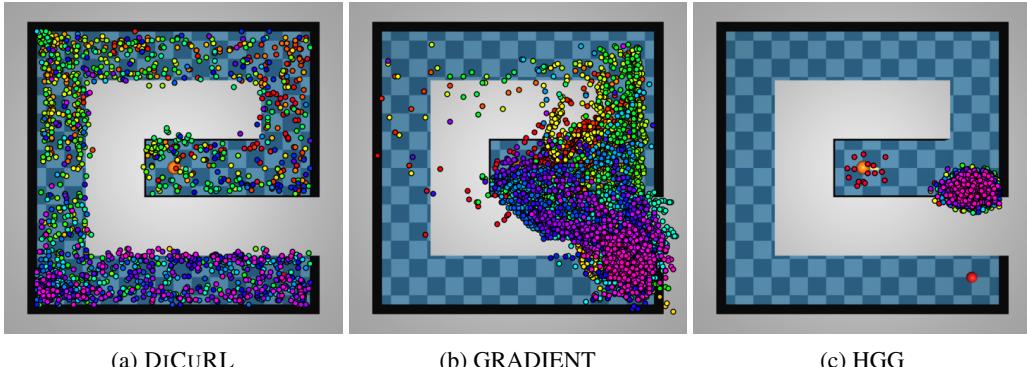


Figure 3: Curriculum goals generated by DiCURL, GRADIENT, and HGG in the PointSpiralMaze environment. The colors ranging from red to purple indicate the curriculum goals across different episodes of the training, and the orange dot and red dot are the agent and desired goal, respectively.

290 Supp. Mat. D shows the curriculum goals generated in the other two environments and illustrates the  
291 dynamics of the diffusion process during training for the case of PointUMaze.

## 292 6 Conclusion and Limitations

293 In this work, we introduced DiCURL, a novel approach that utilizes diffusion models to generate  
294 curriculum goals for an RL agent. The diffusion model is trained to minimize its loss function while  
295 simultaneously maximizing the expected value of  $Q$  and the AIM reward function. The minimization  
296 of the diffusion loss helps capture the visited state distribution, while the maximization of  $Q$  and  
297 AIM reward function helps generate goals at an appropriate difficulty level and at the same time  
298 guide the generated curriculum goals closer to the desired goal. Furthermore, the generated goals  
299 promote exploration due to the inherent noising and denoising mechanism of the diffusion model.  
300 Our proposed approach has two main limitations. First, while diffusion models excel at handling high-  
301 dimensional data, such as images [43], incorporating the AIM reward into a combined loss function  
302 can hinder the curriculum goal generation in such settings as the AIM reward may underperform  
303 in higher dimensionalities. Secondly, we employ the Minimum Cost Maximum Flow algorithm  
304 to select the optimal curriculum goals from the set generated by the diffusion model. However,  
305 alternative selection strategies could potentially be more effective. Future work will aim to address  
306 these limitations and extend our approach to more complex environments.

307 **References**

- 308 [1] Mnih, Volodymyr and Kavukcuoglu, Koray and Silver, David and Graves, Alex and Antonoglou,  
309 Ioannis and Wierstra, Daan and Riedmiller, Martin. Playing Atari with deep reinforcement  
310 learning, 2013. arXiv:1312.5602.
- 311 [2] Mnih, Volodymyr and Kavukcuoglu, Koray and Silver, David and Rusu, Andrei A and Veness,  
312 Joel and Bellemare, Marc G and Graves, Alex and Riedmiller, Martin and Fidjeland, Andreas K  
313 and Ostrovski, Georg and others. Human-level control through deep reinforcement learning.  
314 *Nature*, 518:529–533, 2015.
- 315 [3] Silver, David and Hubert, Thomas and Schrittwieser, Julian and Antonoglou, Ioannis and Lai,  
316 Matthew and Guez, Arthur and Lanctot, Marc and Sifre, Laurent and Kumaran, Dharshan and  
317 Graepel, Thore and others. A general reinforcement learning algorithm that masters chess,  
318 shogi, and Go through self-play. *Science*, 362(6419):1140–1144, 2018.
- 319 [4] Silver, David and Huang, Aja and Maddison, Chris J and Guez, Arthur and Sifre, Laurent and  
320 Van Den Driessche, George and Schrittwieser, Julian and Antonoglou, Ioannis and Panneershel-  
321 vam, Veda and Lanctot, Marc and others. Mastering the game of Go with deep neural networks  
322 and tree search. *Nature*, 529:484–489, 2016.
- 323 [5] Rajeswaran, Aravind and Lowrey, Kendall and Todorov, Emanuel V and Kakade, Sham M.  
324 Towards generalization and simplicity in continuous control. In *Advances in Neural Information  
325 Processing Systems*, volume 30, San Diego, CA, USA, 2017. Neural Information Processing  
326 Systems Foundation.
- 327 [6] Ng, Andrew Y and Coates, Adam and Diel, Mark and Ganapathi, Varun and Schulte, Jamie  
328 and Tse, Ben and Berger, Eric and Liang, Eric. Autonomous inverted helicopter flight via  
329 reinforcement learning. In *International Symposium on Experimental Robotics*, pages 363–372,  
330 Berlin Heidelberg, Germany, 2006. Springer.
- 331 [7] Timothy P. Lillicrap and Jonathan J. Hunt and Alexander Pritzel and Nicolas Heess and Tom  
332 Erez and Yuval Tassa and David Silver and Daan Wierstra. Continuous control with deep  
333 reinforcement learning, 2019. arXiv:1509.02971.
- 334 [8] Sayar, Erdi and Bing, Zhenshan and D’Eramo, Carlo and Oguz, Ozgur S and Knoll, Alois.  
335 Contact Energy Based Hindsight Experience Prioritization, 2023. arXiv:2312.02677.
- 336 [9] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transfor-  
337 mations: Theory and application to reward shaping. In *Icml*, volume 99, pages 278–287. Citeseer,  
338 1999.
- 339 [10] Rui Zhao and Volker Tresp. Energy-based hindsight experience prioritization. In Aude Billard,  
340 Anca Dragan, Jan Peters, and Jun Morimoto, editors, *Proceedings of The 2nd Conference on  
341 Robot Learning*, volume 87 of *Proceedings of Machine Learning Research*, pages 113–122.  
342 PMLR, 29–31 Oct 2018.
- 343 [11] Rui Zhao, Xudong Sun, and Volker Tresp. Maximum entropy-regularized multi-goal rein-  
344 force-ment learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the  
345 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine  
346 Learning Research*, pages 7553–7562. PMLR, 09–15 Jun 2019.
- 347 [12] Daesol Cho, Seungjae Lee, and H. Jin Kim. Outcome-directed reinforcement learning by  
348 uncertainty & temporal distance-aware curriculum goal generation. In *The Eleventh Interna-  
349 tional Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.*  
350 OpenReview.net, 2023.
- 351 [13] Peide Huang, Mengdi Xu, Jiacheng Zhu, Laixi Shi, Fei Fang, and DING ZHAO. Curriculum  
352 reinforcement learning using optimal transport via gradual domain adaptation. In S. Koyejo,  
353 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural  
354 Information Processing Systems*, volume 35, pages 10656–10670. Curran Associates, Inc.,  
355 2022.

- 356 [14] Carlos Florensa, David Held, Xinyang Geng, and Pieter Abbeel. Automatic goal generation  
 357 for reinforcement learning agents. In *International conference on machine learning*, pages  
 358 1515–1528. PMLR, 2018.
- 359 [15] Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In Marina  
 360 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine  
 361 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4940–4950. PMLR,  
 362 18–24 Jul 2021.
- 363 [16] Yunzhi Zhang, Pieter Abbeel, and Lerrel Pinto. Automatic curriculum learning through value  
 364 disagreement. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors,  
 365 *Advances in Neural Information Processing Systems*, volume 33, pages 7648–7659. Curran  
 366 Associates, Inc., 2020.
- 367 [17] Rémy Portelas, Cédric Colas, Katja Hofmann, and Pierre-Yves Oudeyer. Teacher algorithms  
 368 for curriculum learning of deep rl in continuously parameterized environments. In Leslie Pack  
 369 Kaelbling, Danica Kragic, and Komei Sugiura, editors, *Proceedings of the Conference on Robot  
 370 Learning*, volume 100 of *Proceedings of Machine Learning Research*, pages 835–853. PMLR,  
 371 30 Oct–01 Nov 2020.
- 372 [18] Alex Graves, Marc G. Bellemare, Jacob Menick, Rémi Munos, and Koray Kavukcuoglu.  
 373 Automated curriculum learning for neural networks. In Doina Precup and Yee Whye Teh,  
 374 editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of  
 375 *Proceedings of Machine Learning Research*, pages 1311–1320. PMLR, 06–11 Aug 2017.
- 376 [19] Ren, Zhizhou and Dong, Kefan and Zhou, Yuan and Liu, Qiang and Peng, Jian. Exploration via  
 377 Hindsight Goal Generation. arXiv:1906.04279.
- 378 [20] Pascal Klink, Hany Abdulsamad, Boris Belousov, Carlo D'Eramo, Jan Peters, and Joni Pajarinen.  
 379 A probabilistic interpretation of self-paced learning with applications to reinforcement learning.  
 380 *Journal of Machine Learning Research*, 22(182):1–52, 2021.
- 381 [21] Pascal Klink, Carlo D' Eramo, Jan R Peters, and Joni Pajarinen. Self-paced deep reinforcement  
 382 learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances  
 383 in Neural Information Processing Systems*, volume 33, pages 9216–9227. Curran Associates,  
 384 Inc., 2020.
- 385 [22] Pascal Klink, Hany Abdulsamad, Boris Belousov, and Jan Peters. Self-paced contextual  
 386 reinforcement learning. In Leslie Pack Kaelbling, Danica Kragic, and Komei Sugiura, editors,  
 387 *Proceedings of the Conference on Robot Learning*, volume 100 of *Proceedings of Machine  
 388 Learning Research*, pages 513–529. PMLR, 30 Oct–01 Nov 2020.
- 389 [23] Pascal Klink, Haoyi Yang, Carlo D'Eramo, Jan Peters, and Joni Pajarinen. Curriculum rein-  
 390 forcement learning via constrained optimal transport. In Kamalika Chaudhuri, Stefanie Jegelka,  
 391 Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th Inter-  
 392 national Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning  
 393 Research*, pages 11341–11358. PMLR, 17–23 Jul 2022.
- 394 [24] Kevin Li, Abhishek Gupta, Ashwin Reddy, Vitchyr H Pong, Aurick Zhou, Justin Yu, and Sergey  
 395 Levine. Mural: Meta-learning uncertainty-aware rewards for outcome-driven reinforcement  
 396 learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International  
 397 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,  
 398 pages 6346–6356. PMLR, 18–24 Jul 2021.
- 399 [25] Meng Fang, Tianyi Zhou, Yali Du, Lei Han, and Zhengyou Zhang. Curriculum-guided hindsight  
 400 experience replay. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and  
 401 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran  
 402 Associates, Inc., 2019.
- 403 [26] Zhenshan Bing, Matthias Brucker, Fabrice O. Morin, Rui Li, Xiaojie Su, Kai Huang, and  
 404 Alois Knoll. Complex robotic manipulation via graph-based hindsight goal generation. *IEEE  
 405 Transactions on Neural Networks and Learning Systems*, 33(12):7863–7876, 2022.

- 406 [27] Zhenshan Bing, Hongkuan Zhou, Rui Li, Xiaojie Su, Fabrice O. Morin, Kai Huang, and Alois  
 407 Knoll. Solving robotic manipulation with sparse reward reinforcement learning via graph-based  
 408 diversity and proximity. *IEEE Transactions on Industrial Electronics*, 70(3):2759–2769, 2023.
- 409 [28] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion  
 410 for flexible behavior synthesis. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba  
 411 Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Confer-  
 412 ence on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages  
 413 9902–9915. PMLR, 17–23 Jul 2022.
- 414 [29] Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive  
 415 policy class for offline reinforcement learning. *arXiv preprint arXiv:2208.06193*, 2022.
- 416 [30] Chang Chen, Fei Deng, Kenji Kawaguchi, Caglar Gulcehre, and Sungjin Ahn. Simple hierarchi-  
 417 cal planning with diffusion. *arXiv preprint arXiv:2401.02644*, 2024.
- 418 [31] Ishan Durugkar, Mauricio Tec, Scott Niekum, and Peter Stone. Adversarial intrinsic motivation  
 419 for reinforcement learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and  
 420 J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34,  
 421 pages 8622–8636. Curran Associates, Inc., 2021.
- 422 [32] Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based  
 423 control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages  
 424 5026–5033. IEEE, 2012.
- 425 [33] Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E. Taylor, and Peter  
 426 Stone. Curriculum learning for reinforcement learning domains: A framework and survey.  
 427 *CoRR*, abs/2003.04960, 2020.
- 428 [34] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,  
 429 Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience  
 430 replay. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and  
 431 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran  
 432 Associates, Inc., 2017.
- 433 [35] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
 434 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani,  
 435 M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural  
 436 Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- 437 [36] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function ap-  
 438 proximators. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International  
 439 Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*,  
 440 pages 1312–1320, Lille, France, 07–09 Jul 2015. PMLR.
- 441 [37] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adap-  
 442 tation of deep networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the  
 443 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine  
 444 Learning Research*, pages 1126–1135. PMLR, 06–11 Aug 2017.
- 445 [38] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
 446 policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy  
 447 and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine  
 448 Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1861–1870. PMLR,  
 449 10–15 Jul 2018.
- 450 [39] Villani, Cédric and others. *Optimal transport: old and new*, volume 338. Springer, Berlin  
 451 Heidelberg, Germany, 2009.
- 452 [40] Arjovsky, Martin and Chintala, Soumith and Bottou, Léon. Wasserstein generative adversarial  
 453 networks. In *International Conference on Machine Learning*, pages 214–223, Sydney, Australia,  
 454 2017. PMLR.

- 455 [41] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In  
456 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural*  
457 *Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- 458 [42] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for  
459 statisticians. *Journal of the American Statistical Association*, 112(518):859–877, April 2017.
- 460 [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
461 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*  
462 *conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 463 [44] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsu-  
464 pervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei,  
465 editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of  
466 *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France, 07–09 Jul 2015.  
467 PMLR.
- 468 [45] Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint*  
469 *arXiv:2208.11970*, 2022.
- 470 [46] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
471 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

472 **Supplemental Material**

473 **A Conditional Diffusion Model**

474 Diffusion models [41, 44] are usually defined as Markov chains and trained using variational inference.  
 475 The forward diffusion process gradually adds Gaussian noise to the data, while the reverse diffusion  
 476 process subtracts the learned noise from the data.

477 Given the goal data  $\mathbf{g}_0$  sampled from the real goal distribution  $\mathbf{g}_0 \sim \mathcal{G}$ , we can define the forward  
 478 diffusion process using Markov chains and add a small amount of Gaussian noise in each timestep  
 479  $k$  according to a variance schedule  $\{\beta_k \in (0, 1)\}_{k=1}^N$ , and generate a sequence of noisy goal data  
 480 samples  $\mathbf{g}_1, \dots, \mathbf{g}_N$ :

$$\begin{aligned} q(\mathbf{g}_{1:N} | \mathbf{g}_0) &:= \prod_{k=1}^N q(\mathbf{g}_k | \mathbf{g}_{k-1}) \\ q(\mathbf{g}_k | \mathbf{g}_{k-1}) &:= \mathcal{N}\left(\mathbf{g}_k; \sqrt{1 - \beta_k} \mathbf{g}_{k-1}, \beta_k \mathbf{I}\right). \end{aligned} \quad (13)$$

- 481 As the index of the timestep  $k$  increases, the goal sample  $\mathbf{g}_0$  gradually loses its features and becomes  
 482 an isotropic Gaussian noise when  $N \rightarrow \infty$ .  
 483 Instead of sampling data recursively for a given timestep  $k$ , we can sample  $\mathbf{g}_k$  in a closed form using  
 484 the reparameterization trick<sup>5</sup>:

$$\mathbf{g}_k \sim \mathcal{N}\left(\mathbf{g}_k; \sqrt{1 - \beta_k} \mathbf{g}_{k-1}, \beta_k \mathbf{I}\right). \quad (14)$$

485 Let  $\alpha_k = 1 - \beta_k$  and  $\bar{\alpha}_k = \prod_{i=1}^N \alpha_i$ . We can write:

$$\mathbf{g}_k = \sqrt{\alpha_k} \mathbf{g}_{k-1} + \sqrt{1 - \alpha_k} \epsilon_{k-1} \quad (15)$$

$$\mathbf{g}_k = \sqrt{\alpha_k \alpha_{k-1}} \mathbf{g}_{k-2} + \sqrt{1 - \alpha_k} \epsilon_{k-1} + \sqrt{\alpha_k (1 - \alpha_{k-1})} \epsilon_{k-2}. \quad (16)$$

487 We then merge two Gaussian distributions with different variance  $\mathcal{N}(0, \sigma_1^2)$  and  $\mathcal{N}(0, \sigma_2^2)$  into a new  
 488 Gaussian distribution  $\mathcal{N}(0, \sigma_1^2 + \sigma_2^2)$ . Specifically, we can merge the two equations above, where  
 489  $\sqrt{1 - \alpha_k} + \alpha_k (1 - \alpha_{k-1}) = \sqrt{1 - \alpha_k \alpha_{k-1}}$ :

$$\mathbf{g}_k = \sqrt{\alpha_k \alpha_{k-1}} \mathbf{g}_{k-2} + \sqrt{1 - \alpha_k \alpha_{k-1}} \bar{\epsilon}_{k-2} \quad (17)$$

$$\mathbf{g}_k = \sqrt{\prod_{i=1}^N \alpha_i} \mathbf{g}_0 + \sqrt{1 - \prod_{i=1}^N \alpha_i} \epsilon_0 \quad (18)$$

491 where  $\epsilon_k \sim \mathcal{N}(0, 1)$ .

492 Now, we can write:

$$\mathbf{g}_k = \sqrt{\bar{\alpha}_k} \mathbf{g}_0 + \sqrt{1 - \bar{\alpha}_k} \epsilon \quad (19)$$

$$q(\mathbf{g}_k | \mathbf{g}_0) = \mathcal{N}(\mathbf{g}_k; \sqrt{\bar{\alpha}_k} \mathbf{g}_0, (1 - \bar{\alpha}_k) \mathbf{I}) \quad (20)$$

494 Using the equation above, we can sample data in a closed form at any arbitrary timestep  $k$ .

495 As mentioned in the main text, the reverse diffusion process aims to recover the original input data  
 496 from the noisy (diffused) data. It learns to progressively reverse the diffusion process step by step,  
 497 and approximates the joint distribution  $p_\theta(\mathbf{x}_{0:T})$ . This process is defined as:

$$\begin{aligned} p_\theta(\mathbf{g}_{0:N}) &:= p(\mathbf{g}_N) \prod_{k=1}^T p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k) \\ p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k) &:= \mathcal{N}(\mathbf{g}_{k-1}; \boldsymbol{\mu}_\theta(\mathbf{g}_k, k), \boldsymbol{\Sigma}_\theta(\mathbf{g}_k, k)) \end{aligned} \quad (21)$$

---

<sup>5</sup>The reparameterization trick works by separating the deterministic and the stochastic parts of the sampling operation. Instead of directly sampling from the distribution  $z \sim \mathcal{N}(\mu, \sigma)$ , we can sample  $\epsilon$  from the normal distribution  $\mathcal{N}(0, 1)$ , multiply it by standard deviation  $\sigma$ , and add mean  $\mu$ :  $z = \mu + \sigma \epsilon$

498 The reverse conditional probability is tractable when conditioned by  $g^0$ .

499 By Bayes rule, we have:

$$q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0) = \frac{q(\mathbf{g}_k | \mathbf{g}_{k-1}, \mathbf{g}_0)q(\mathbf{g}_{k-1} | \mathbf{g}_0)}{q(\mathbf{g}_k | \mathbf{g}_0)}. \quad (22)$$

500 As we assume that the diffusion model is a Markov chain, the future state depends only on the present  
501 state, namely:

$$q(\mathbf{g}_k | \mathbf{g}_{k-1}, \mathbf{g}_0) = q(\mathbf{g}_k | \mathbf{g}_{k-1}) = \mathcal{N}(\mathbf{g}_k; \sqrt{\alpha_k} \mathbf{g}_{k-1}, (1 - \alpha_k) \mathbf{I}) \quad (23)$$

$$q(\mathbf{g}_{k-1} | \mathbf{g}_0) = \mathcal{N}(\mathbf{g}_{k-1}; \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0, (1 - \bar{\alpha}_{k-1}) \mathbf{I}) \quad (24)$$

$$q(\mathbf{g}_k | \mathbf{g}_0) = \mathcal{N}(\mathbf{g}_k; \sqrt{\bar{\alpha}_k} \mathbf{g}_0, (1 - \bar{\alpha}_k) \mathbf{I}) \quad (25)$$

502

$$q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0) = \frac{\mathcal{N}(\mathbf{g}_k; \sqrt{\alpha_k} \mathbf{g}_{k-1}, (1 - \alpha_k) \mathbf{I}) \mathcal{N}(\mathbf{g}_{k-1}; \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0, (1 - \bar{\alpha}_{k-1}) \mathbf{I})}{\mathcal{N}(\mathbf{g}_k; \sqrt{\bar{\alpha}_k} \mathbf{g}_0, (1 - \bar{\alpha}_k) \mathbf{I})} \quad (26)$$

503 where:

$$\mathcal{N}(\mathbf{g}_k; \sqrt{\alpha_k} \mathbf{g}_{k-1}, (1 - \alpha_k) \mathbf{I}) = \frac{1}{(2\pi)^{n/2} |(1 - \alpha_k) \mathbf{I}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{g}_k - \sqrt{\alpha_k} \mathbf{g}_{k-1})^T \frac{\mathbf{I}}{(1 - \alpha_k)} (\mathbf{g}_k - \sqrt{\alpha_k} \mathbf{g}_{k-1})\right) \quad (27)$$

$$\approx \exp\left(-\frac{(\mathbf{g}_k - \sqrt{\alpha_k} \mathbf{g}_{k-1})^2}{2(1 - \alpha_k) \mathbf{I}}\right) \quad (28)$$

$$\mathcal{N}(\mathbf{g}_{k-1}; \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0, (1 - \bar{\alpha}_{k-1}) \mathbf{I}) = \frac{1}{(2\pi)^{n/2} |(1 - \bar{\alpha}_{k-1}) \mathbf{I}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{g}_{k-1} - \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0)^T \frac{\mathbf{I}}{(1 - \bar{\alpha}_{k-1})} (\mathbf{g}_{k-1} - \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0)\right) \quad (29)$$

$$\approx \exp\left(-\frac{(\mathbf{g}_{k-1} - \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0)^2}{2(1 - \bar{\alpha}_{k-1})}\right) \quad (30)$$

$$\mathcal{N}(\mathbf{g}_k; \sqrt{\bar{\alpha}_k} \mathbf{g}_0, (1 - \bar{\alpha}_k) \mathbf{I}) = \frac{1}{(2\pi)^{n/2} |(1 - \bar{\alpha}_k) \mathbf{I}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{g}_k - \sqrt{\bar{\alpha}_k} \mathbf{g}_0)^T \frac{\mathbf{I}}{(1 - \bar{\alpha}_k)} (\mathbf{g}_k - \sqrt{\bar{\alpha}_k} \mathbf{g}_0)\right) \quad (31)$$

$$\approx \exp\left(-\frac{(\mathbf{g}_k - \sqrt{\bar{\alpha}_k} \mathbf{g}_0)^2}{2(1 - \bar{\alpha}_k)}\right) \quad (32)$$

504 We then can write:

$$q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0) \approx \exp\left(-\frac{1}{2\mathbf{I}} \left( \frac{(\mathbf{g}_k - \sqrt{\alpha_k} \mathbf{g}_{k-1})_2}{(1 - \alpha_k)} + \frac{(\mathbf{g}_{k-1} - \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0)_2}{1 - \bar{\alpha}_{k-1}} - \frac{(\mathbf{g}_k - \sqrt{\bar{\alpha}_k} \mathbf{g}_0)_2}{1 - \bar{\alpha}_k} \right)\right) \quad (33)$$

$$\approx \exp\left(-\frac{1}{2\mathbf{I}} \left( \frac{(\mathbf{g}_k - \sqrt{\alpha_k} \mathbf{g}_{k-1})^2}{(1 - \alpha_k)} + \frac{(\mathbf{g}_{k-1} - \sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0)^2}{1 - \bar{\alpha}_{k-1}} - \frac{(\mathbf{g}_k - \sqrt{\bar{\alpha}_k} \mathbf{g}_0)^2}{1 - \bar{\alpha}_k} \right)\right) \quad (34)$$

$$\approx \exp\left(-\frac{1}{2\mathbf{I}} \left( \frac{(\mathbf{g}_k^2 - 2\sqrt{\alpha_k} \mathbf{g}_k \mathbf{g}_{k-1} + \alpha_k \mathbf{g}_{k-1}^2)}{(1 - \alpha_k)} + \frac{(\mathbf{g}_{k-1}^2 - 2\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_{k-1} \mathbf{g}_0 + \bar{\alpha}_{k-1} \mathbf{g}_0^2)}{1 - \bar{\alpha}_{k-1}} - \frac{(\mathbf{g}_k - \sqrt{\bar{\alpha}_k} \mathbf{g}_0)^2}{1 - \bar{\alpha}_k} \right)\right) \quad (35)$$

$$= \exp\left(-\frac{1}{2\mathbf{I}} \left( \frac{(-2\sqrt{\alpha_k} \mathbf{g}_k \mathbf{g}_{k-1} + \alpha_k \mathbf{g}_{k-1}^2)}{(1 - \alpha_k)} + \frac{(\mathbf{g}_{k-1}^2 - 2\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_{k-1} \mathbf{g}_0)}{1 - \bar{\alpha}_{k-1}} + \mathcal{C}(\mathbf{g}_k, \mathbf{g}_0) \right)\right) \quad (36)$$

$$\approx \exp\left(-\frac{1}{2\mathbf{I}} \left( -\frac{2\sqrt{\alpha_k} \mathbf{g}_k \mathbf{g}_{k-1}}{(1 - \alpha_k)} + \frac{\alpha_k \mathbf{g}_{k-1}^2}{(1 - \alpha_k)} + \frac{\mathbf{g}_{k-1}^2}{1 - \bar{\alpha}_{k-1}} - \frac{2\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_{k-1} \mathbf{g}_0}{1 - \bar{\alpha}_{k-1}} \right)\right) \quad (37)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{\alpha_k}{1-\alpha_k} + \frac{1}{1-\bar{\alpha}_{k-1}} \right) \mathbf{g}_{k-1}^2 - 2 \left( \frac{\sqrt{\alpha_k} \mathbf{g}_k}{1-\alpha_k} + \frac{\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0}{1-\bar{\alpha}_{k-1}} \right) \mathbf{g}_{k-1} \right) \right) \quad (38)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{\alpha_k(1-\bar{\alpha}_{k-1})+1-\alpha_k}{(1-\alpha_k)(1-\bar{\alpha}_{k-1})} \right) \mathbf{g}_{k-1}^2 - 2 \left( \frac{\sqrt{\alpha_k} \mathbf{g}_k}{1-\alpha_k} + \frac{\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0}{1-\bar{\alpha}_{k-1}} \right) \mathbf{g}_{k-1} \right) \right) \quad (39)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{\alpha_k-\bar{\alpha}_k+1-\alpha_k}{(1-\alpha_k)(1-\bar{\alpha}_{k-1})} \right) \mathbf{g}_{k-1}^2 - 2 \left( \frac{\sqrt{\alpha_k} \mathbf{g}_k}{1-\alpha_k} + \frac{\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0}{1-\bar{\alpha}_{k-1}} \right) \mathbf{g}_{k-1} \right) \right) \quad (40)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{1-\bar{\alpha}_k}{(1-\alpha_k)(1-\bar{\alpha}_{k-1})} \right) \mathbf{g}_{k-1}^2 - 2 \left( \frac{\sqrt{\alpha_k} \mathbf{g}_k}{1-\alpha_k} + \frac{\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0}{1-\bar{\alpha}_{k-1}} \right) \mathbf{g}_{k-1} \right) \right) \quad (41)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{1-\bar{\alpha}_k}{(1-\alpha_k)(1-\bar{\alpha}_{k-1})} \right) \left( \mathbf{g}_{k-1}^2 - 2 \left( \frac{\frac{\sqrt{\alpha_k} \mathbf{g}_k}{1-\alpha_k} + \frac{\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0}{1-\bar{\alpha}_{k-1}}}{\frac{1-\bar{\alpha}_k}{(1-\alpha_k)(1-\bar{\alpha}_{k-1})}} \right) \mathbf{g}_{k-1} \right) \right) \right) \quad (42)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{1-\bar{\alpha}_k}{(1-\alpha_k)(1-\bar{\alpha}_{k-1})} \right) \left( \mathbf{g}_{k-1}^2 - 2 \frac{\left( \frac{\sqrt{\alpha_k} \mathbf{g}_k}{1-\alpha_k} + \frac{\sqrt{\bar{\alpha}_{k-1}} \mathbf{g}_0}{1-\bar{\alpha}_{k-1}} \right) (1-\alpha_k)(1-\bar{\alpha}_{k-1})}{1-\bar{\alpha}_k} \mathbf{g}_{k-1} \right) \right) \right) \quad (43)$$

$$= \exp \left( -\frac{1}{2\mathbf{I}} \left( \left( \frac{1}{\frac{(1-\alpha_k)(1-\bar{\alpha}_{k-1})}{1-\bar{\alpha}_k}} \right) \left( \mathbf{g}_{k-1}^2 - 2 \frac{\sqrt{\alpha_k}(1-\bar{\alpha}_{k-1})\mathbf{g}_k + \sqrt{\bar{\alpha}_{k-1}}(1-\alpha_k)\mathbf{g}_0}{1-\bar{\alpha}_k} \mathbf{g}_{k-1} \right) \right) \right) \quad (44)$$

$$\approx \mathcal{N} \left( \mathbf{g}_{k-1}; \underbrace{\frac{\sqrt{\alpha_k}(1-\bar{\alpha}_{k-1})\mathbf{g}_k + \sqrt{\bar{\alpha}_{k-1}}(1-\alpha_k)\mathbf{g}_0}{1-\bar{\alpha}_k}}_{\mu(\mathbf{g}_k, \mathbf{g}_0)}, \underbrace{\frac{(1-\alpha_k)(1-\bar{\alpha}_{k-1})}{1-\bar{\alpha}_k} \mathbf{I}}_{\Sigma(k)} \right) \quad (45)$$

505 where  $C(\mathbf{g}_k, \mathbf{g}_0)$  is a constant term with respect to  $\mathbf{g}_{k-1}$  and is calculated using only  $\mathbf{g}_k$ ,  $\mathbf{g}_0$ , and  $\alpha_k$ .

506 The Eq. (45) shows that at every step  $\mathbf{g}_{k-1} \sim q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0)$  follows a normal distribution with  
507 mean  $\mu(\mathbf{g}_k, \mathbf{g}_0)$  and variance  $\Sigma(k)$ .

508 Given this, the goal is to learn a denoising model  $p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k)$  by approximating the ground-truth  
509 denoising transition  $q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0)$ . This can be done by minimizing the KL divergence between  
510  $p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k)$  and  $q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0)$ .

511 The KL divergence between two Gaussian distributions  $\mathcal{N}(x; \mu_x, \Sigma_x)$  and  $\mathcal{N}(y; \mu_y, \Sigma_y)$  is defined  
512 as:

$$D_{\text{KL}}(\mathcal{N}(x; \mu_x, \Sigma_x) || \mathcal{N}(y; \mu_y, \Sigma_y)) = \frac{1}{2} \left[ \log \frac{|\Sigma_y|}{|\Sigma_x|} - d + \text{tr}(\Sigma_y^{-1} \Sigma_x) + (\mu_y - \mu_x)^T \Sigma_y^{-1} (\mu_y - \mu_x) \right]. \quad (46)$$

513 The KL divergence between the denoising model  $p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k)$  and the ground-truth denoising  
514 transition  $q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0)$  is:

$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0) || p_\theta(\mathbf{g}_{k-1} | \mathbf{g}_k)). \quad (47)$$

515 The learning of the denoising model can be written as:

$$\arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mathbf{g}_{k-1}; \mu(\mathbf{g}_k, \mathbf{g}_0), \Sigma) || \mathcal{N}(\mathbf{g}_{k-1}; \mu_\theta, \Sigma)) \quad (48)$$

516

$$\arg \min_{\theta} \frac{1}{2} \left[ \log \frac{|\Sigma|}{|\Sigma|} - d + \text{tr}(\Sigma^{-1} \Sigma) + (\mu_\theta - \mu(\mathbf{g}_k, \mathbf{g}_0))^T \Sigma^{-1} (\mu_\theta - \mu(\mathbf{g}_k, \mathbf{g}_0)) \right] \quad (49)$$

517

$$\arg \min_{\theta} \frac{1}{2} [\log 1 - d + d + (\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0))^T \Sigma^{-1} (\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0))] \quad (50)$$

518

$$\arg \min_{\theta} \frac{1}{2} [(\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0))^T \Sigma^{-1} (\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0))] \quad (51)$$

519

$$\arg \min_{\theta} \frac{1}{2} [(\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0))^T (\mathbf{I})^{-1} (\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0))] \quad (52)$$

520

$$\arg \min_{\theta} \frac{1}{2} [\|\mu_{\theta} - \mu(\mathbf{g}_k, \mathbf{g}_0)\|_2^2] . \quad (53)$$

521 We can rearrange Eq. (19) as follows:

$$\mathbf{g}_0 = \frac{\mathbf{g}_k - \sqrt{1 - \bar{\alpha}_k} \epsilon}{\sqrt{\bar{\alpha}_k}} \quad (54)$$

522 and rewrite  $\mu(\mathbf{g}_k, \mathbf{g}_0)$  in Eq. (45) as follows:

$$\mu(\mathbf{g}_k, \mathbf{g}_0) = \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})\mathbf{g}_k + \sqrt{\bar{\alpha}_{k-1}}(1 - \alpha_k) \frac{\mathbf{g}_k - \sqrt{1 - \bar{\alpha}_k} \epsilon}{\sqrt{\bar{\alpha}_k}}}{1 - \bar{\alpha}_k} \quad (55)$$

523 where  $\bar{\alpha}_k = \prod_{i=1}^N \alpha_i$ .

524 Then:

$$\mu(\mathbf{g}_k, \mathbf{g}_0) = \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})\mathbf{g}_k + \sqrt{\alpha_1 \alpha_2 \dots \alpha_{k-1}}(1 - \alpha_k) \frac{\mathbf{g}_k - \sqrt{1 - \bar{\alpha}_k} \epsilon}{\sqrt{\alpha_1 \alpha_2 \dots \alpha_k}}}{1 - \bar{\alpha}_k} \quad (56)$$

$$= \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})\mathbf{g}_k}{1 - \bar{\alpha}_k} + \frac{(1 - \alpha_k) \mathbf{g}_k}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} - \frac{(1 - \alpha_k) \sqrt{1 - \bar{\alpha}_k} \epsilon}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} \quad (57)$$

$$= \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})\mathbf{g}_k}{1 - \bar{\alpha}_k} + \frac{(1 - \alpha_k) \mathbf{g}_k}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} - \frac{(1 - \alpha_k) \sqrt{1 - \bar{\alpha}_k} \epsilon}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} \quad (58)$$

$$= \left( \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k} + \frac{(1 - \alpha_k)}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} \right) \mathbf{g}_k - \frac{(1 - \alpha_k) \sqrt{1 - \bar{\alpha}_k} \epsilon}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} \quad (59)$$

$$= \left( \frac{\alpha_k(1 - \bar{\alpha}_{k-1})}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} + \frac{1 - \alpha_k}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} \right) \mathbf{g}_k - \frac{(1 - \alpha_k)}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon \quad (60)$$

$$= \left( \frac{\alpha_k - \bar{\alpha}_k + 1 - \alpha_k}{(1 - \bar{\alpha}_k) \sqrt{\alpha_k}} \right) \mathbf{g}_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon \quad (61)$$

$$= \frac{1}{\sqrt{\alpha_k}} \mathbf{g}_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon \quad (62)$$

525 We can model the denoising transition mean as follows:

$$\mu_{\theta}(\mathbf{g}_k, k) = \frac{1}{\sqrt{\alpha_k}} \mathbf{g}_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon_{\theta}(\mathbf{g}_k, k) . \quad (63)$$

526 Given this, the minimization of the KL divergence can be written as:

527

$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{g}_{k-1} | \mathbf{g}_k, \mathbf{g}_0) || p_{\theta}(\mathbf{g}_{k-1} | \mathbf{g}_k)) \quad (64)$$

528

$$\arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mathbf{g}_{k-1}; \mu(\mathbf{g}_k, \mathbf{g}_0), \Sigma) || \mathcal{N}(\mathbf{g}_{k-1}; \mu_{\theta}, \Sigma)) \quad (65)$$

529

$$\arg \min_{\theta} \frac{1}{2} \left[ \left\| \frac{1}{\sqrt{\alpha_k}} \mathbf{g}_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon_{\theta}(\mathbf{g}_k, k) - \frac{1}{\sqrt{\alpha_k}} \mathbf{g}_k + \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon \right\|_2^2 \right] \quad (66)$$

$$\arg \min_{\theta} \frac{1}{2} \left[ \left\| \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon_{\theta}(\mathbf{g}_k, k) \right\|_2^2 \right] \quad (67)$$

530

$$\arg \min_{\theta} \frac{1}{2} \frac{(1 - \alpha_k)^2}{(1 - \bar{\alpha}_k) \alpha_k} \left[ \|\epsilon - \epsilon_{\theta}(\mathbf{g}_k, k)\|_2^2 \right]. \quad (68)$$

531 Thus far, we have only modeled the goal distribution  $p_{\theta}(\mathbf{g})$ . To transform this into a conditional  
 532 diffusion model, we can incorporate the state  $s$  information at each diffusion timestep  $k$  [45] in Eq.  
 533 (21):

$$p_{\theta}(\mathbf{g}_{0:N}) := p(\mathbf{g}_N) \prod_{k=1}^T p_{\theta}(\mathbf{g}_{k-1} | \mathbf{g}_k s) \quad (69)$$

$$p_{\theta}(\mathbf{g}_{k-1} | \mathbf{g}_k) := \mathcal{N}(\mathbf{g}_{k-1}; \boldsymbol{\mu}_{\theta}(\mathbf{g}_k, s, k), \boldsymbol{\Sigma}_{\theta}(\mathbf{g}_k, s, k)).$$

534 To conclude, we can model the denoising transition mean by rewriting Eq. (70) based on the state  $s$   
 535 as follows:

$$\boldsymbol{\mu}_{\theta}(\mathbf{g}_k, s, k) = \frac{1}{\sqrt{\alpha_k}} \mathbf{g}_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon_{\theta}(\mathbf{g}_k, s, k). \quad (70)$$

## 536 B $s_T$ Strategy

537 In contrast to selecting optimal curriculum goals via bipartite graph optimization, as shown in Eq. 12  
 538 in the main text<sup>6</sup> we performed an additional experiment where we considered only the state from the  
 539 final timestep  $T$ , under the assumption that it is closer to achieving the desired goal. The state  $s_T$  at  
 540 the last timestep of each episode is then input into the diffusion model, which generates curriculum  
 541 goals to be achieved in the subsequent episode. The resulting variant is detailed in Algorithm 3, with  
 542 the RL training process described in Algorithm 4.

543 As demonstrated by this additional experiment, the assumption that the final state is closer to the  
 544 desired goal does not always hold true. As a matter of fact, the state at the last timestep in some  
 545 episodes might be even further away from the desired goal compared to the states in previous timesteps.  
 546 Consequently, the curriculum goals generated by the diffusion model may not gradually shift from  
 547 the initial position toward the desired goal. If the state at the last timestep is not progressively moving  
 548 towards the desired goal, the curriculum goals generated by the diffusion model may jump to different  
 549 areas within the maze environment. This can lead to a decrease in sample efficiency and a slower  
 550 success rate. This issue is illustrated in Fig. 4, where it can be seen, that particularly in the case  
 551 of PointSpiralMaze, the variant of our algorithm that makes use of this strategy (that we call “ $s_T$   
 552 strategy”) converges later than the main described in the main text.

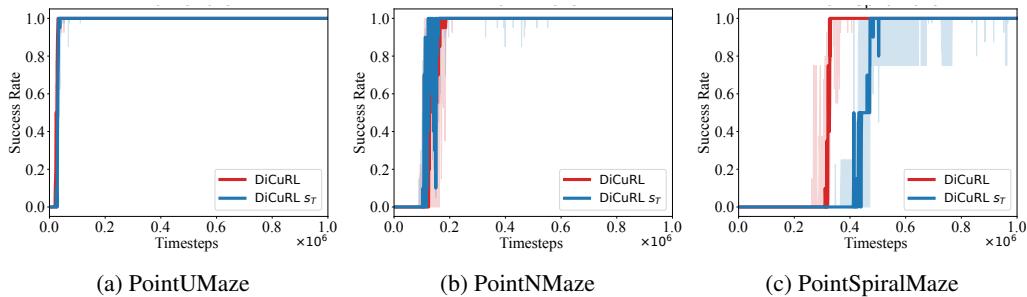


Figure 4: Test success rate for the original DiCuRL algorithm described in the main text compared to a variant that makes use of the  $s_T$  strategy.

<sup>6</sup>As a reminder, in the proposed algorithm we sample a mini-batch  $b$  (containing many states from different timesteps) from the replay buffer  $\mathcal{B}$  and provide it to the curriculum goal generator (i.e., the diffusion model) to generate a curriculum goal distribution  $\mathcal{G}_c$ , and then select the optimal curriculum goal  $g_c$  using bipartite graph optimization by maximizing the cost function given in Eq. 12.

---

**Algorithm 3** Diffusion Curriculum Goal Generator with the  $s_T$  Strategy

---

```
1: Input: no. of reverse diffusion timestep  $N$ , training step  $M$ 
2: for  $i = 1, \dots, M$  do ▷ Training iterations
3:    $\mathbf{g}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:   for  $k = N, \dots, 1$  do ▷ Reverse diffusion process
5:      $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
6:      $g_{k-1} = \frac{1}{\sqrt{\alpha_k}} \left( g_k - \frac{\beta_k}{\sqrt{1-\alpha_k}} \epsilon_\theta(g_k, k, \textcolor{brown}{s}_T) \right) + \sqrt{\beta_k} \epsilon$  ▷ using Eq. (9)
7:      $k \sim \text{Uniform}(\{1, \dots, N\})$ 
8:      $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
9:   Calculate diffusion  $\mathcal{L}_d(\theta), Q(s, \mathbf{g}_0, \pi(s)), r(\mathbf{g}_0, g_d)$  ▷ Calculate with the generated goal  $\mathbf{g}_0$ 
10:  Calculate total loss  $\mathcal{L} = \xi_d \mathcal{L}_d(\theta) - \xi_q Q(s, \mathbf{g}_0, \pi(s)) - \xi_r r(\mathbf{g}_0, g_d)$  ▷ Eq 11
11:   $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}$  ▷ Take gradient descent step and update the diffusion weights
12: return  $\mathbf{g}_0$ 
```

---

---

**Algorithm 4** RL Training with the  $s_T$  Strategy

---

```
1: Input: no. of episodes  $E$ , timesteps  $T$ 
2: Select an off-policy algorithm  $\mathbb{A}$  ▷ In our case,  $\mathbb{A}$  is SAC
3: Initialize replay buffer  $\mathcal{B} \leftarrow \emptyset$ ,  $g_c \leftarrow g_d$  and networks  $Q_\phi, \pi_\psi, r_\varphi$ 
4: for episode  $= 0 \dots E$  do
5:   Sample initial state  $s_0$ 
6:   for  $t = 0 \dots T$  do
7:      $a_t = \pi(s_t, g_c)$ 
8:     Execute  $a_t$ , obtain next state  $s_{t+1}$ 
9:     Store transition  $(s_t, a_t, r_t, s_{t+1}, g_c)$  in  $\mathcal{B}$ 
10:    Sample a minibatch  $b$  from replay buffer  $\mathcal{B}$ 
11:     $g_c \leftarrow \text{DiffusionCurriculumGenerator}(\textcolor{brown}{s}_T)$ 
12:    Update  $Q$  and  $\pi$  with  $b$  to minimize  $\mathcal{L}_Q$  and  $\mathcal{L}_\pi$  in Eq. (1) and in Eq. (2)
13:    Update the AIM reward function  $r_\varphi$ 
14:    $success \leftarrow 0$  ▷ Success Rate
15:   Sample a desired goal  $g_d \sim \mathcal{G}$ 
16:   for  $i = 1 \dots n_{testrollout}$  do
17:      $a_t = \pi(s_t, g_d)$ 
18:     Execute  $a_t$ , obtain next state  $s_{t+1}$  and reward  $r_t$ 
19:     if  $|\phi(s_{t+1}) - g_d| \leq \kappa$  then
20:        $success = success + 1/n_{testrollout}$ 
```

---

553 **C Experimental Details**

554 For implementing DICURL, we utilized the original implementation of OUTPACE [12] and aug-  
555 mented this codebase with our diffusion model for curriculum goal generation. We conducted our  
556 experiments on a cluster computer using an NVIDIA RTX A5000 GPU, 64GB of RAM, and a 4-core  
557 CPU. For the PointSpiralMaze environment, the total compute time of the main method and the  $s_T$   
558 strategy is approximately 23 hours and 5 hours, respectively. For the remaining environments, the  
559 main method and the  $s_T$  strategy require approximately 11.5 hours and 2.5 hours, respectively.

560 **Baselines.** The baseline CRL algorithms are trained as follows:

- 561 • HGG [19]: We utilized the default settings from the original implementation, which can be found  
562 at <https://github.com/Stilwell-Git/Hindsight-Goal-Generation>.
- 563 • CURROT [23]: We adhered to the default settings of the original implementation, available at  
564 <https://github.com/psclklnk/currot>.
- 565 • ALP-GMM [17], VDS [16], PLR [15], ACL [18], GoalGAN [14]: We followed the default settings  
566 from the implementation available at <https://github.com/psclklnk/currot>.

- 567 • GRADIENT [13]: We used the default settings from the implementation available at <https://github.com/PeideHuang/gradient>.  
 568  
 569 • OUTPACE [12]: We implemented the default settings from the implementation available at  
 570 [https://github.com/jayLEE0301/outpace\\_official](https://github.com/jayLEE0301/outpace_official).

571 **Training details.** All baseline models, except for the GRADIENT method, were trained using the  
 572 Soft Actor-Critic (SAC) [38]. Although the original implementations of these algorithms primarily  
 573 use on-policy methods, adaptations were made in the CURROT repository to utilize the off-policy  
 574 SAC algorithm. This modification, available in the CURROT repository, allows for a more direct  
 575 comparison of sample efficiency across all models. The GRADIENT method was trained using  
 576 both the SAC and Proximal Policy Optimization (PPO) algorithms [46]. However, we present only  
 577 the results obtained with the PPO algorithm, as it outperforms SAC in the three considered maze  
 578 environments. All our algorithms' hyperparameters used in the experiments are reported in Table 3.

Table 3: Hyperparameters for DiCURL.

Parameter	Value	Parameter	Value
Critic hidden dim	512	Discount factor $\gamma_r$	0.99
Critic hidden depth	3	$r_\varphi$ Update frequency	1000
Critic target $\tau$	0.01	No. of gradient steps for $r_\varphi$ update	10
Critic target update frequency	2	No. of ensemble networks for $r_\varphi$	5
Actor hidden dim	512	learning rate for $r_\varphi$	1e-4
Actor hidden depth	3	RL optimizer	Adam
Actor update frequency	2	Diffusion network updates per episode	25
RL batch size	512	Diffusion loss coefficient $\xi_d$	1
Init. temperature $\alpha_{\text{init}}$ of SAC	0.3	$Q$ -function coefficient $\xi_q$	10
Replay buffer $\mathcal{B}$ size	1e6	AIM reward function coefficient $\xi_r$	1
Diffusion training iterations	300	Reverse diffusion timesteps	10
Diffusion learning rate	3e-4	Diffusion loss type	$l_2$
Diffusion update frequency	2500	No. of training timesteps	1e6

579 **Environment details.** In each task, the agent's state  $s$  is represented as a vector, comprising its  
 580 position, velocity along the  $x$  and  $y$  axes, orientation angle, and angular velocity around the  $z$ -axis.  
 581 The agent's actions are determined by its velocity and angular velocity. The agent starts each episode  
 582 from an initial state of  $[0, 0]$ . The desired goal distribution is derived by introducing uniform noise  
 583 to the desired goal position. The desired goal positions are set to  $[0, 8]$ ,  $[8, 16]$ , and  $[8, -8]$ , while  
 584 the dimensions of the map are  $12 \times 12$ ,  $12 \times 20$ ,  $20 \times 20$ , respectively for the PointUMaze, the  
 585 PointNMaze, and the PointSpiralMaze task. A task is considered successful when the agent reaches  
 586 the sampled desired goal within a threshold distance of 0.5.

587 **Detailed results.** We report in Table 4 the detailed results in terms of no. of timesteps needed to  
 588 reach a success rate of 0.3 for all the algorithms under comparison. It can be seen that most of the  
 589 algorithms do not reach this rate within the given budget.

## 590 D Generated Curriculum Goals

591 Fig. 5 and Fig. 6 show the differences between the curriculum goals generated by DiCURL,  
 592 GRADIENT, and HGG, respectively in the PointUMaze and PointNMaze environment. As shown  
 593 in the case of PointSpiralMaze reported in the main text, also in these two environments DiCURL  
 594 manages to generate curriculum goals that explore the whole environment more effectively than  
 595 GRADIENT and HDD.

### 596 D.1 Dynamics of the Diffusion Process during Training

597 Considering the case of the PointUMaze environment, the curriculum goal sets, denoted as  
 598  $g_9, g_8, \dots, g_0$  in Fig. 7 generated during the initial phase of training ( $\sim 2000$  timesteps), are obtained

Table 4: No. of timesteps (rounded) to reach success rate 0.3 (average  $\pm$  std. dev. across 5 runs). NaN indicates that a success rate of 0.3 is not reached within the maximum budget of 1e6 timesteps.

Algorithm	PointUMaze	PointNMaze	PointSpiralMaze
DiCURL (Ours)	$25000 \pm 3958$	$99000 \pm 32031$	$302000 \pm 44121$
OUTPACE	$27600 \pm 4176$	$106666 \pm 25603$	$389375 \pm 111661$
HGG	$34750 \pm 14956$	$68000 \pm 15000$	NaN
GRADIENT	$166986 \pm 86132$	$755573 \pm 65967$	NaN
CURROT	$700000 \pm 234520$	NaN	NaN
ALP-GMM	$660000 \pm 205912$	NaN	NaN
GoalGAN	$960000 \pm 290516$	NaN	NaN
PLR	$833333 \pm 253859$	NaN	NaN
VDS	$838461 \pm 294927$	NaN	NaN
ACL	$1000000 \pm 316227$	NaN	NaN

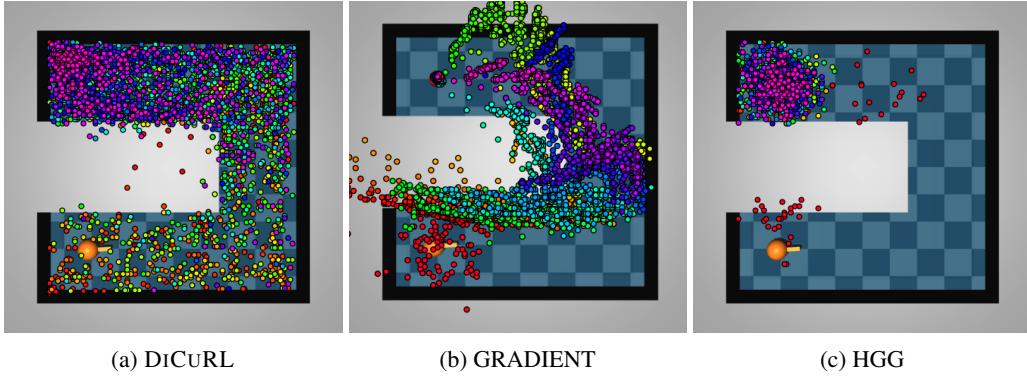


Figure 5: Curriculum goals generated by DiCURL, GRADIENT, and HGG in the PointUMaze environment. The colors ranging from red to purple indicate the curriculum goals across different episodes of the training, and the orange dot and red dot are the agent and desired goal, respectively.

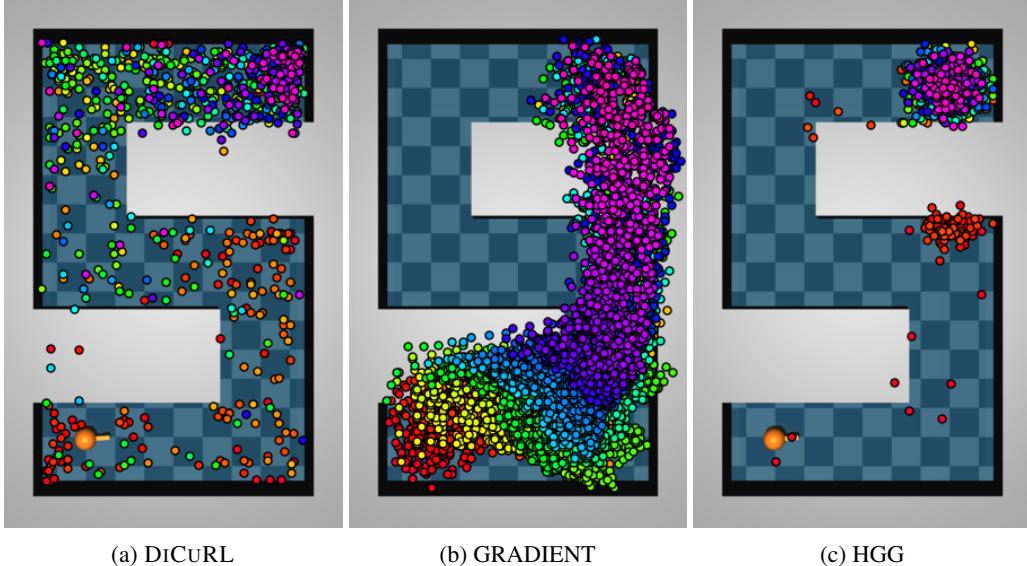


Figure 6: Curriculum goals generated by DiCURL, GRADIENT, and HGG in the PointNMaze environment. The colors ranging from red to purple indicate the curriculum goals across different episodes of the training, and the orange dot and red dot are the agent and desired goal, respectively.

599 through the reverse diffusion process. This process starts from Gaussian noise  $\mathbf{g}_9$  and culminates in  
 600 the final curriculum goal set  $\mathbf{g}_0$ , as expressed in lines 5 and 7 in Algorithm 1. It can be observed that  
 601 the final curriculum goal set  $\mathbf{g}_0$  is generated diagonally. This pattern arises due to the combined loss  
 602 function of the AIM reward and the  $Q$ -value, along with the diffusion loss, as expressed in Eq. (11).  
 603 The AIM reward function, relative to the desired goal  $g_d$ , and the  $Q$ -value, derived from a sampled  
 604 mini-batch  $b$  from the replay buffer  $\mathcal{B}$  at timestep  $\sim 2000$ , are depicted in Fig. 8a and Fig. 8b,  
 605 respectively, using a 10-level contour plot. It can be observed that the reward values increase  
 606 diagonally, while the  $Q$ -value has not yet shown a discernible pattern. As a result, the curriculum  
 607 goal set generated by the diffusion model exhibits a pattern similar to the reward value.

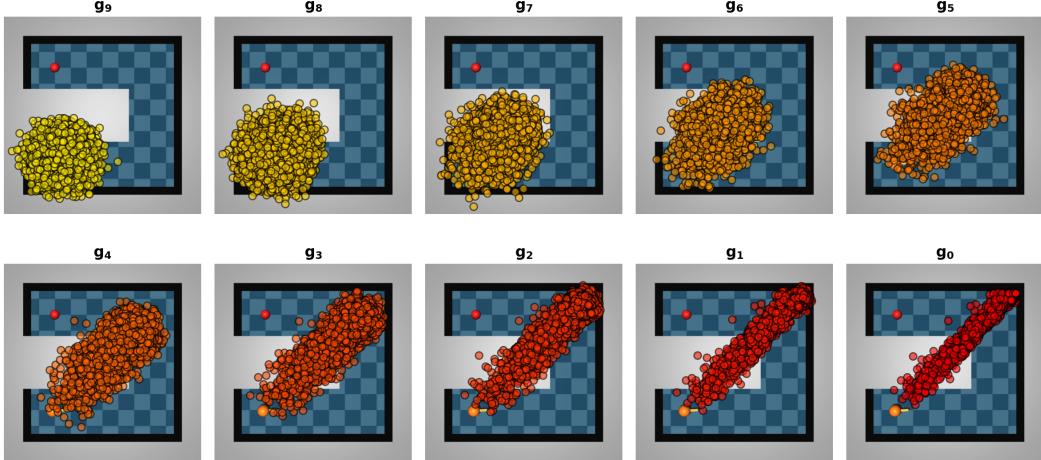


Figure 7: The curriculum goal set  $\mathcal{G}_c$  generated by DiCURL during the reverse diffusion process in the early stage of the training ( $\sim 2000$  timesteps) for the PointUMaze environment.

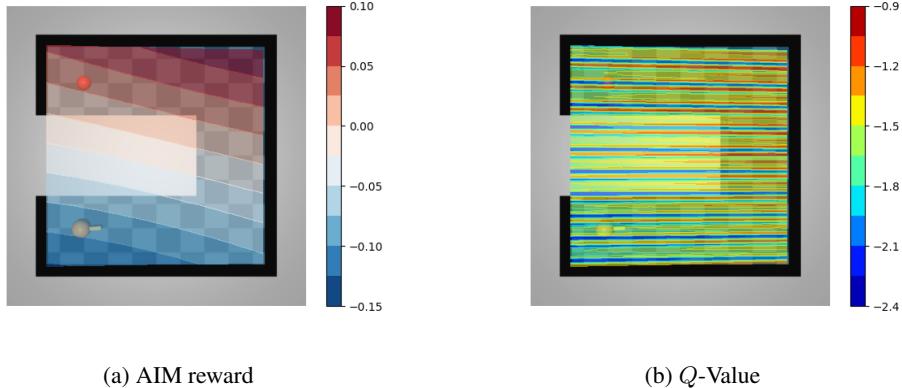


Figure 8: Visualization of the AIM reward function and  $Q$ -function in the early stage of the training ( $\sim 2000$  timesteps) for the PointUMaze environment.

608 In the mid-training stage ( $\sim 15000$  timesteps), the final curriculum goal set  $\mathbf{g}_0$  reflects the explored  
 609 region of the environment, extending up to the top right corner, as shown in Fig. 9. This behavior  
 610 aligns with the agent's exploration progress.  
 611 The AIM reward function in Fig. 10a exhibits higher and better-converged values compared to Fig.  
 612 8a, particularly in the top right corner. Similarly, the  $Q$ -value in Fig. 10b demonstrates its highest  
 613 values concentrated from the middle right side to the top right corner. These observations suggest  
 614 that the final curriculum goal set  $\mathbf{g}_0$  has indeed been formed, encompassing the explored area up to  
 615 the top right corner.

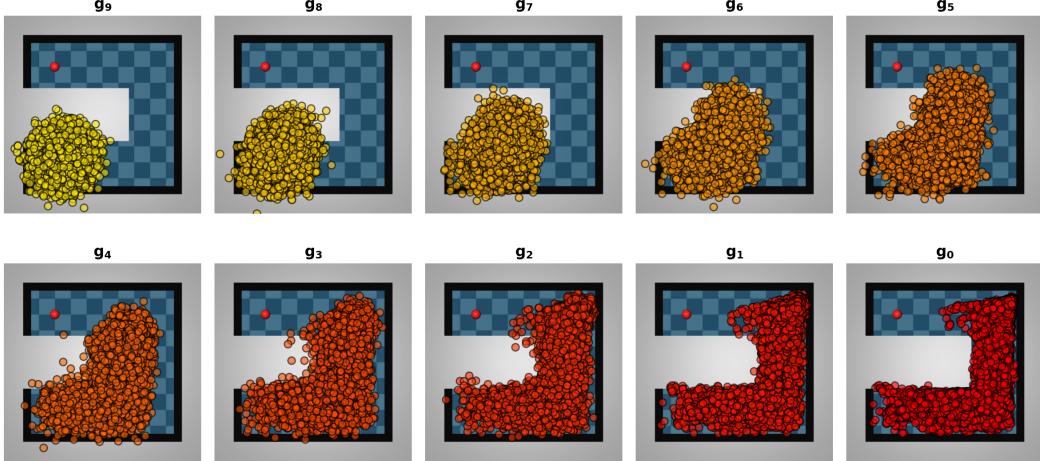


Figure 9: The curriculum goal set  $\mathcal{G}_c$  generated by DiCURL during the reverse diffusion process in the middle stage of the training ( $\sim 15000$  timesteps) for the PointUMaze environment.

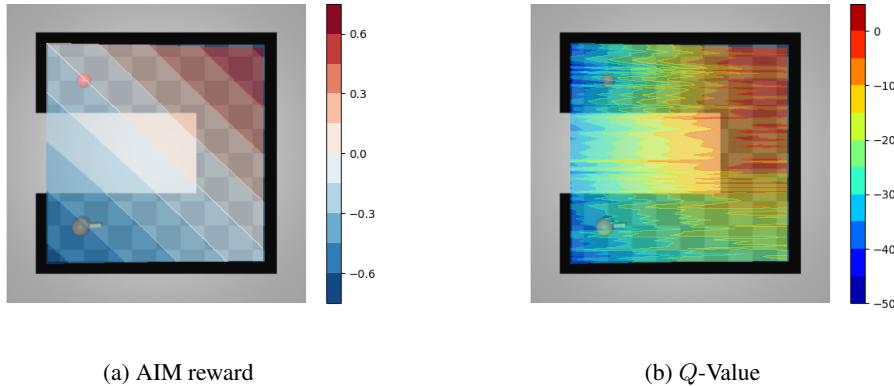


Figure 10: Visualization of the AIM reward function and  $Q$ -function in the middle stage of the training ( $\sim 15000$  timesteps) for the PointUMaze environment.

616 In the near-optimal policy learning phase ( $\sim 30000$  timesteps), the final curriculum goal set  $g_0$   
617 encompasses the entire environment, extending from the initial position to the desired goal, as  
618 depicted in Fig. 11. This behavior aligns with the agent’s progress towards achieving the optimal  
619 policy.

620 The AIM reward function in Fig. 12a appears to be converging to the desired goal more slowly. The  
621  $Q$ -value in Fig. 12b demonstrates its highest values concentrated around the desired goal area. These  
622 observations suggest that the final curriculum goal set  $g_0$  has indeed been formed, encompassing  
623 the entire environment. Furthermore, it appears that the  $Q$ -value converges to the desired goal  
624 more rapidly than the AIM reward function. This observation aligns with our choice of a higher  
625 hyperparameter coefficient  $\xi_q$  for the  $Q$ -value function compared to hyperparameter coefficient  $\xi_r$   
626 for the reward function in the total loss function shown in Eq. (11). This prioritizes the  $Q$ -value’s  
627 influence on shaping the curriculum goals, promoting faster convergence towards the desired behavior.

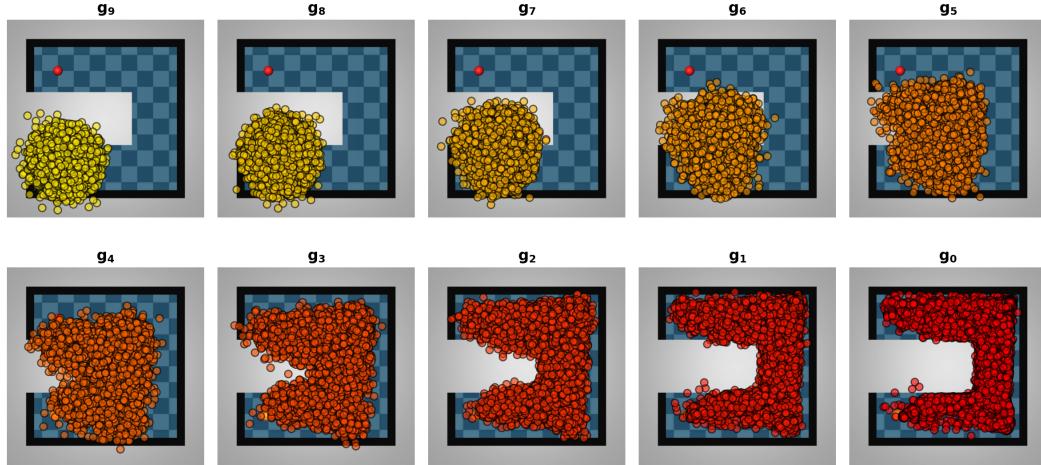


Figure 11: The curriculum goal set  $\mathcal{G}_c$  generated by DiCURL during the reverse diffusion process in the near-optimal stage of the training ( $\sim 30000$  timesteps).

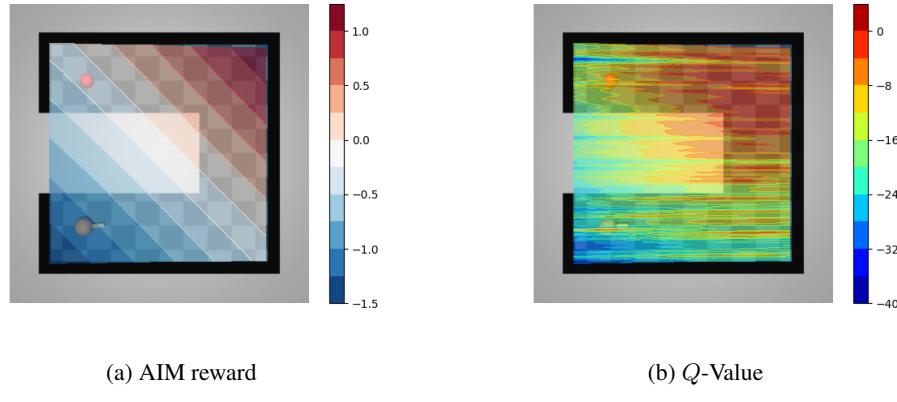


Figure 12: Visualization of the AIM reward function and  $Q$ -function in the near-optimal stage of the training ( $\sim 30000$  timesteps) for the PointUMaze environment.

628 **NeurIPS Paper Checklist**

629 **1. Claims**

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

632 Answer: [Yes]

633 Justification: All claims made in the abstract and introduction accurately reflect the paper's  
634 contributions and scope.

635 **2. Limitations**

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

637 Answer: [Yes]

638 Justification: See Section 6.

639 **3. Theory Assumptions and Proofs**

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

642 Answer: [Yes]

643 Justification: All the theoretical results are detailed in Supp. Mat. A.

644 **4. Experimental Result Reproducibility**

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

648 Answer: [Yes]

649 Justification: We disclose all the experimental details along with our codebase.

650 **5. Open access to data and code**

651 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
652 tions to faithfully reproduce the main experimental results, as described in supplemental  
653 material?

654 Answer: [Yes]

655 Justification: See Supp. Mat. C.

656 **6. Experimental Setting/Details**

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

660 Answer: [Yes]

661 Justification: See Supp. Mat. C.

662 **7. Experiment Statistical Significance**

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

665 Answer: [Yes]

666 Justification: We provide all results in terms of mean and standard deviation.

667 **8. Experiments Compute Resources**

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

671 Answer: [Yes]

672 Justification: See Supp. Mat. C.

673 **9. Code Of Ethics**

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

676 Answer: [Yes]

677 Justification: Our research conforms with the NeurIPS Code of Ethics.

## 678 10. Broader Impacts

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

681 Answer: [Yes]

682 Justification: We summarize the importance and impact of this research in the Introduction.

683 Guidelines:

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

## 706 11. Safeguards

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

710 Answer: [NA]

711 Justification: The paper poses no such risks.

712 Guidelines:

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

## 723 12. Licenses for existing assets

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

727                  Answer: [Yes]

728                  Justification: The creators or original owners of assets (e.g., code, data, models) used in the  
729                  paper are properly credited through appropriate citations or website URLs.

730                  Guidelines:

- 731                  • The answer NA means that the paper does not use existing assets.
- 732                  • The authors should cite the original paper that produced the code package or dataset.
- 733                  • The authors should state which version of the asset is used and, if possible, include a  
734                  URL.
- 735                  • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 736                  • For scraped data from a particular source (e.g., website), the copyright and terms of  
737                  service of that source should be provided.
- 738                  • If assets are released, the license, copyright information, and terms of use in the  
739                  package should be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets)  
740                  has curated licenses for some datasets. Their licensing guide can help determine the  
741                  license of a dataset.
- 742                  • For existing datasets that are re-packaged, both the original license and the license of  
743                  the derived asset (if it has changed) should be provided.
- 744                  • If this information is not available online, the authors are encouraged to reach out to  
745                  the asset's creators.

### 746                  13. New Assets

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

749                  Answer: [Yes]

750                  Justification: We provide our code through an anonymized URL, as well as all the details  
751                  about training and computational setup (see Supp. Mat. C).

752                  Guidelines:

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

### 761                  14. Crowdsourcing and Research with Human Subjects

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

765                  Answer: [NA]

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

767                  Guidelines:

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

### 776                  15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 777                  Subjects

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

782 Answer: [NA]

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

784 Guidelines:

- 785 • The answer NA means that the paper does not involve crowdsourcing nor research with  
786 human subjects.
- 787 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
788 may be required for any human subjects research. If you obtained IRB approval, you  
789 should clearly state this in the paper.
- 790 • We recognize that the procedures for this may vary significantly between institutions  
791 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
792 guidelines for their institution.
- 793 • For initial submissions, do not include any information that would break anonymity (if  
794 applicable), such as the institution conducting the review.