

Diffusion-based Curriculum Reinforcement Learning

Erdi Sayar¹

Giovanni Iacca²

Ozgur S. Oguz³

Alois Knoll¹

¹Technical University of Munich, Germany

²University of Trento, Italy

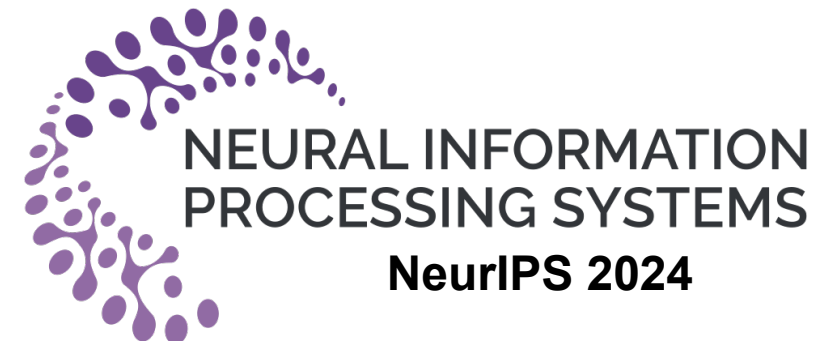
³Bilkent University, Türkiye



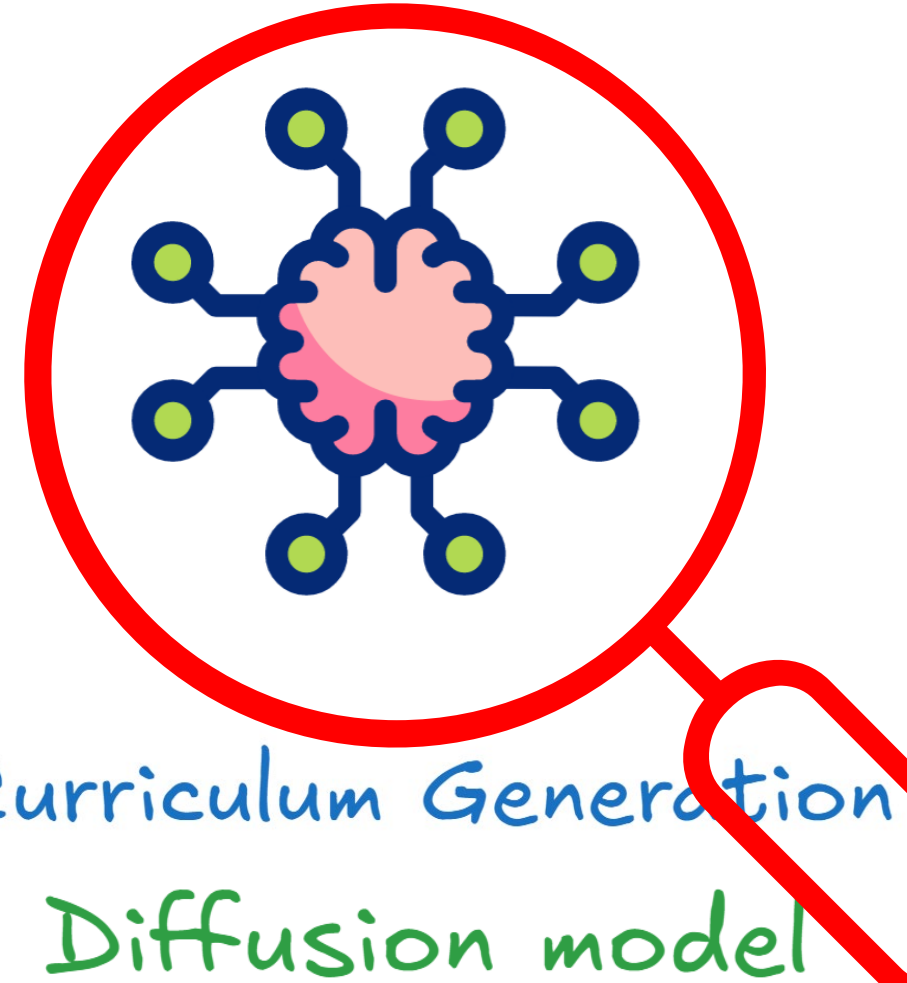
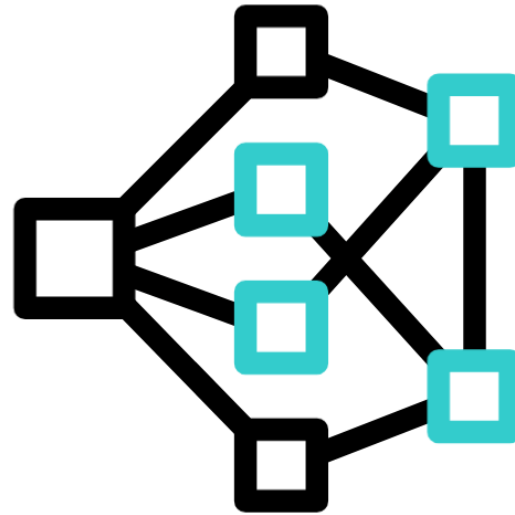
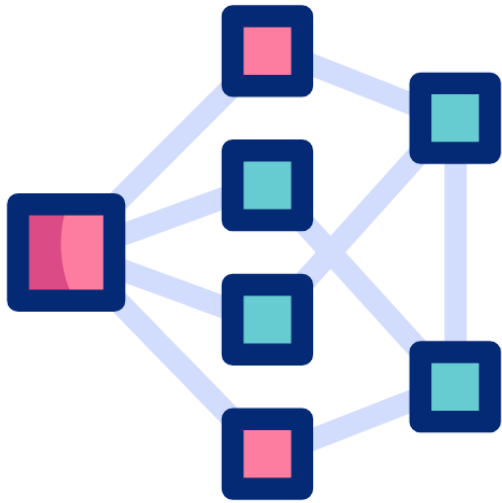
Technische Universität München



UNIVERSITÀ
DI TRENTO



Curriculum generation with Diffusion Models



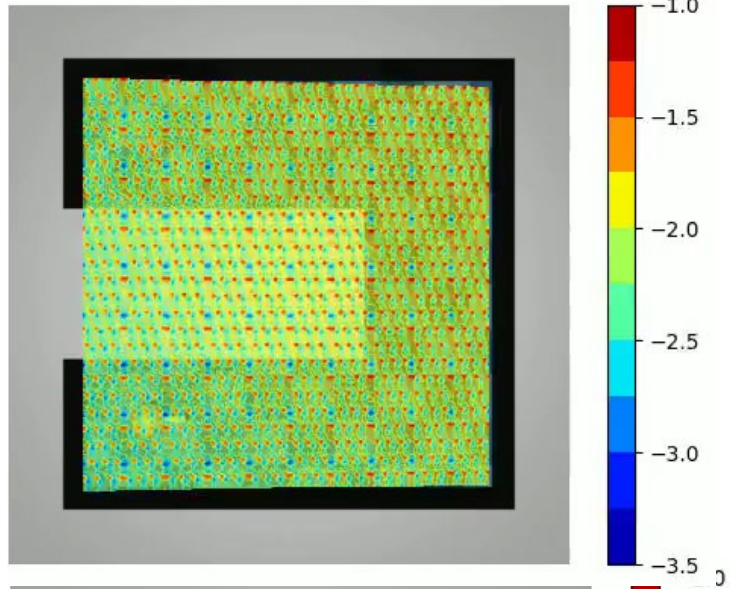
Soft Actor Critic
SAC

AIM Reward
Trainable reward

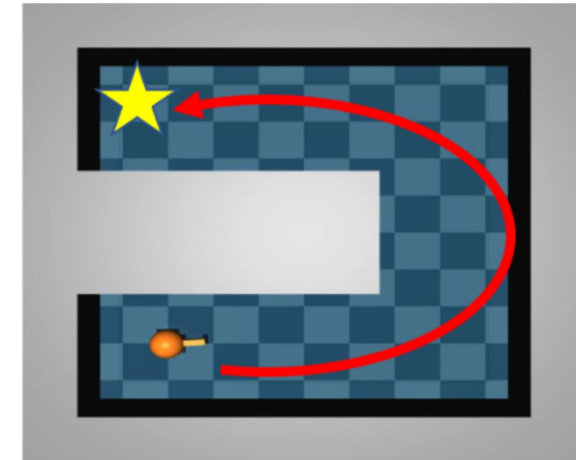
Curriculum Generation
Diffusion model

How Q and reward function changes during training

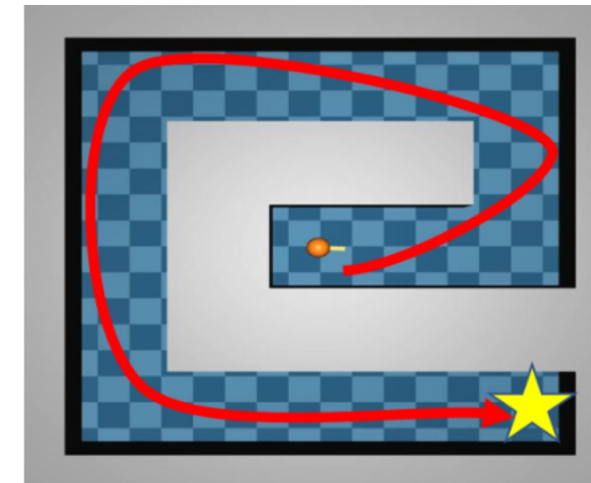
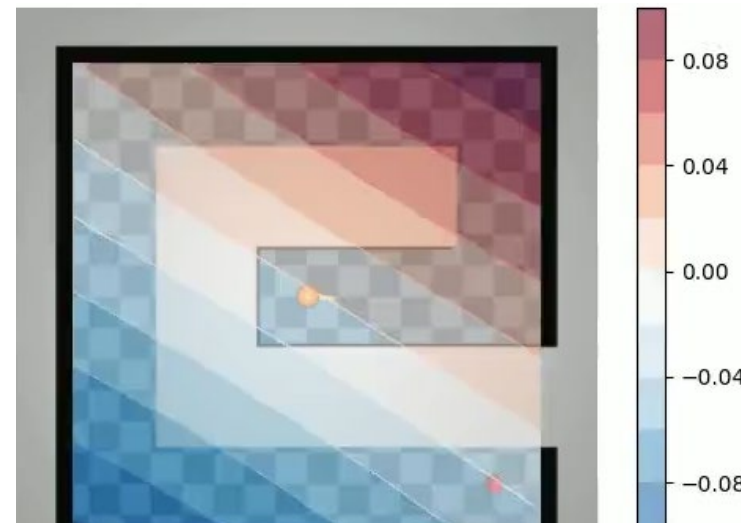
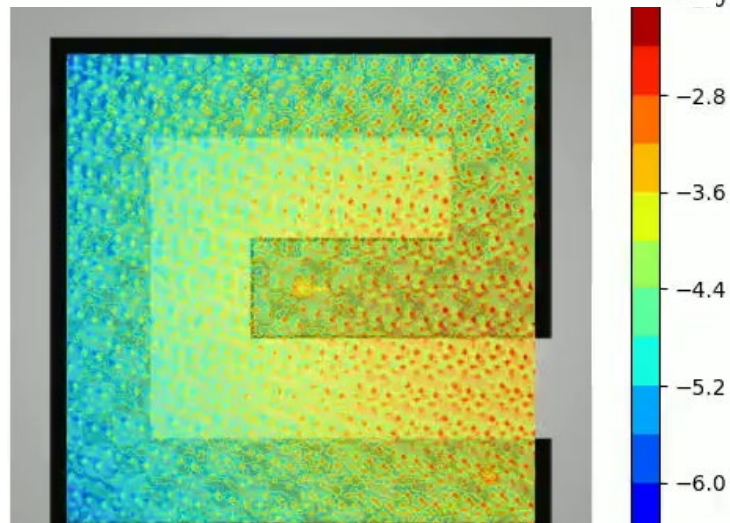
$$Q_{\phi}(s, \mathbf{x}, a)$$



$$r_{\varphi}(\mathbf{x}, g_d)$$



U-Maze Environment



Spiral-Maze Environment

Methodology

- We used a diffusion model to generate curriculum for the agent.

Algorithm 2 RL Training and Evaluation

```
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:         $success = success + 1/n_{testrollout}$ 
```

Methodology

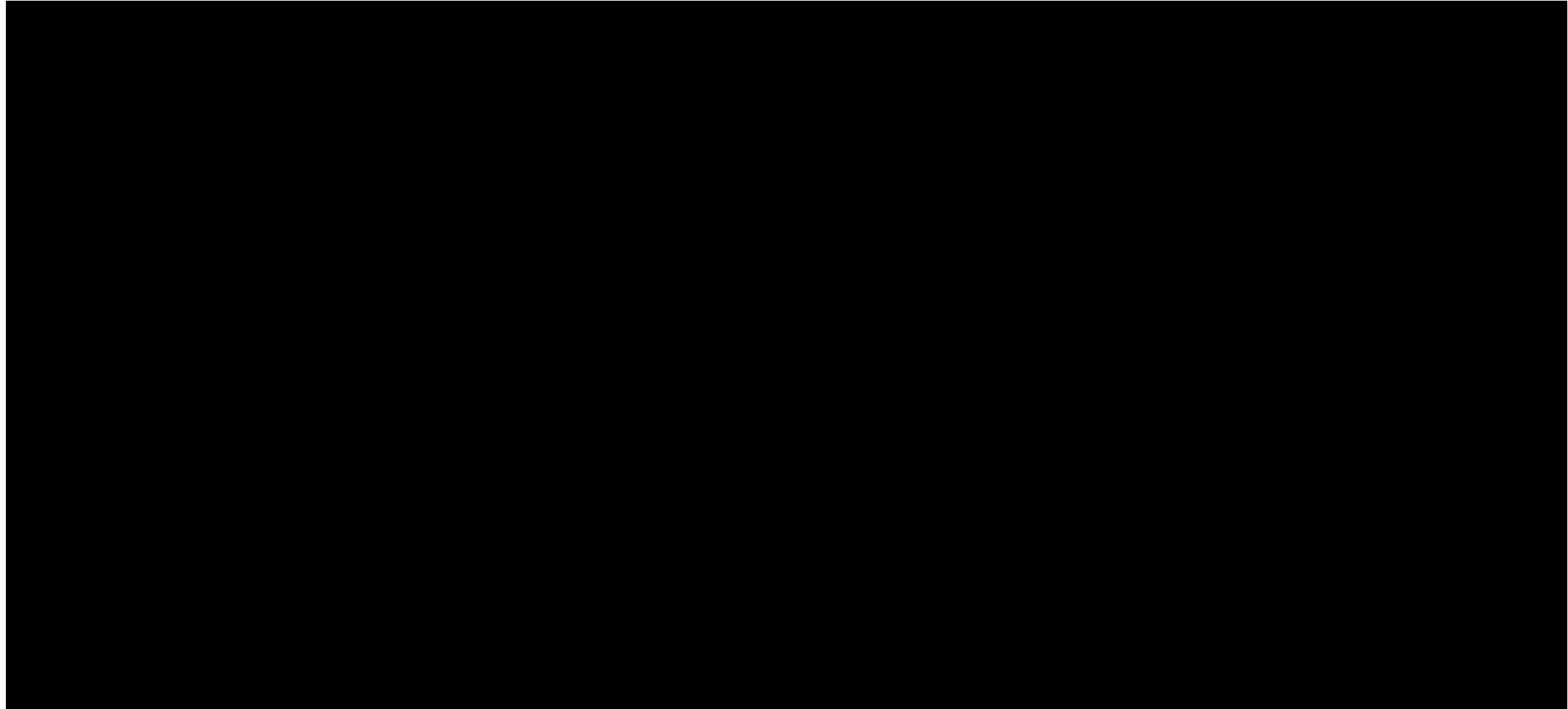
- We used a diffusion model to generate curriculum for the agent.

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  $\mathcal{L}_d(\theta), Q(s, \mathbf{g}_0, \pi(s)), r(g_d, \mathbf{g}_0)$  ▷ 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  $\mathbf{g}_0$ 
```



How Diffusion model is learning to generate curriculum goals.

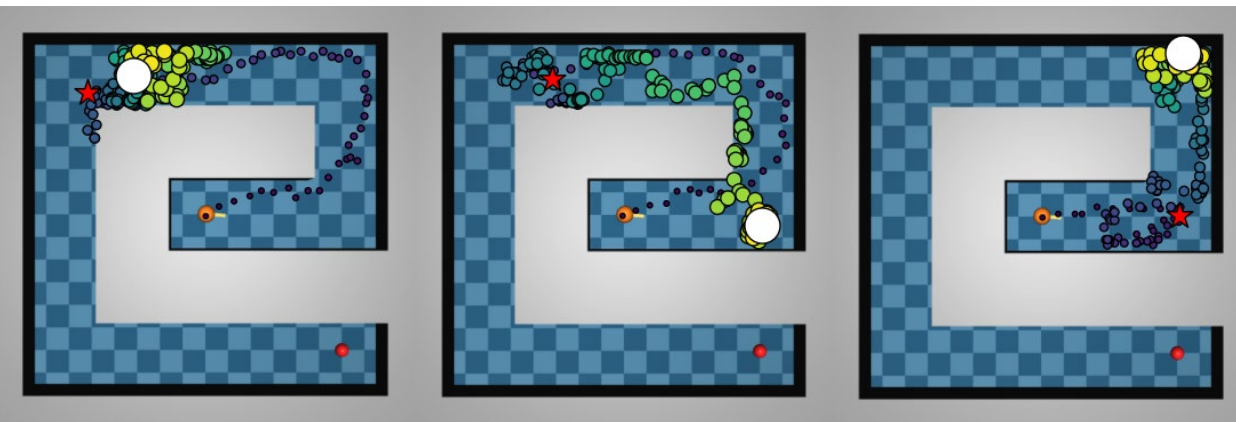


Different Strategies to select curriculum goals

Selecting Curriculum Goals

Strategy s_T

Strategy Optimization Algorithm



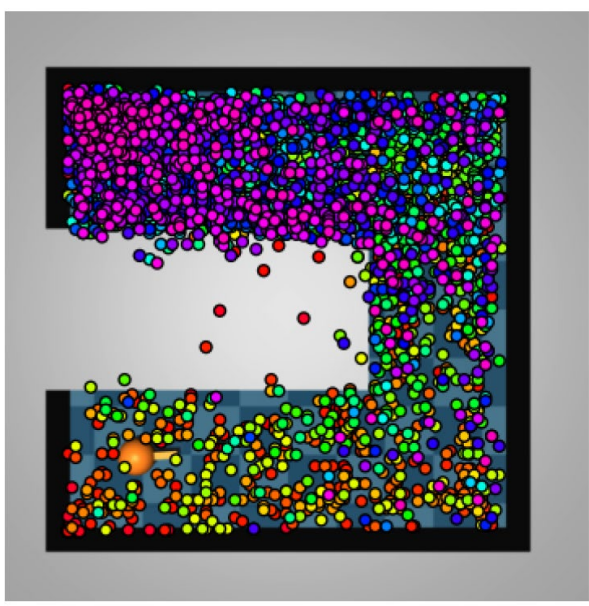
Episode 525

Episode 526

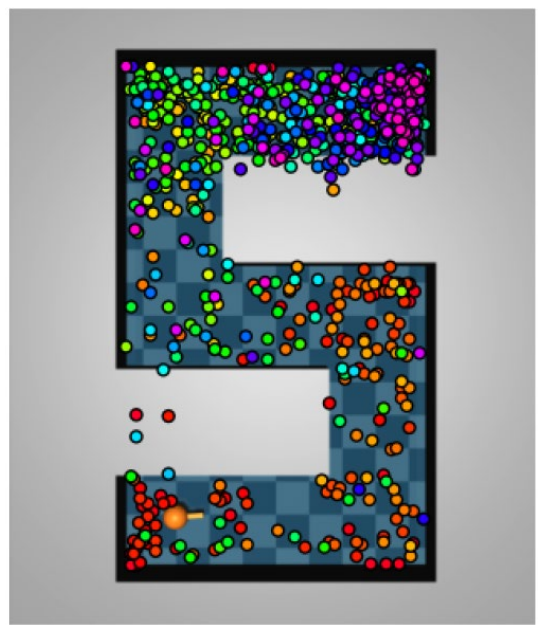
Episode 527

$$\max_{\hat{\mathcal{G}}^c: |\hat{\mathcal{G}}^c|=K} \sum_{g_{t=0, \dots, T}^0 \in \hat{\mathcal{G}}^d, g_+^i \in \hat{\mathcal{G}}^+} w \quad \text{where: } w = \sqrt{\frac{(g_i - \bar{g})^2}{N}}$$

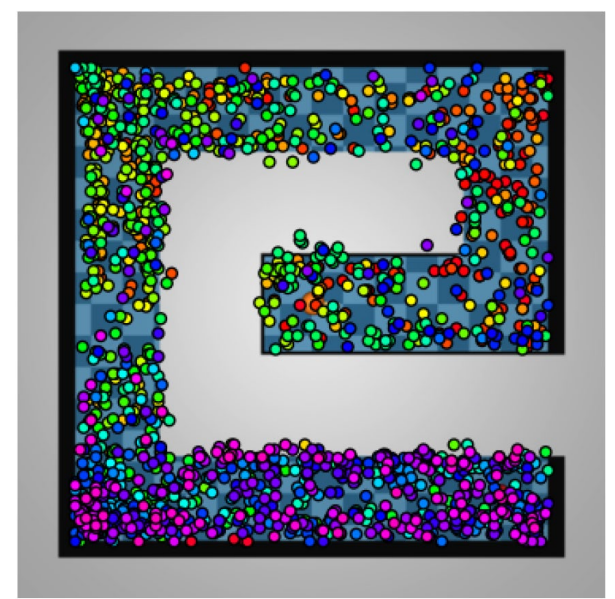
Generated curriculum goals during training



PointUMaze

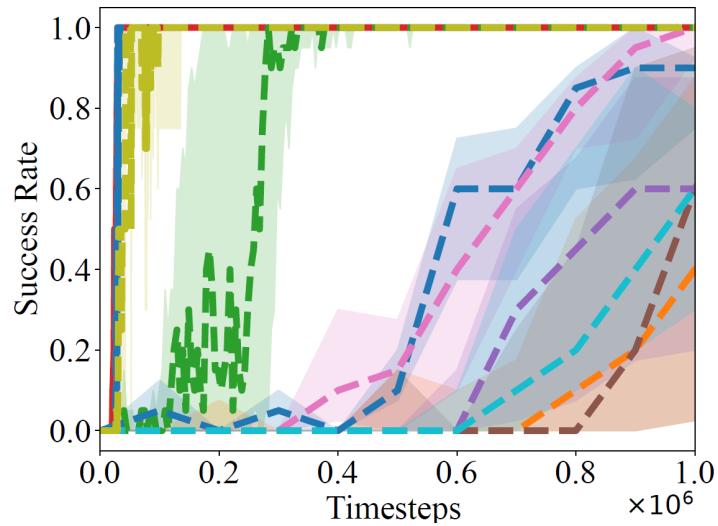


PointNMaze

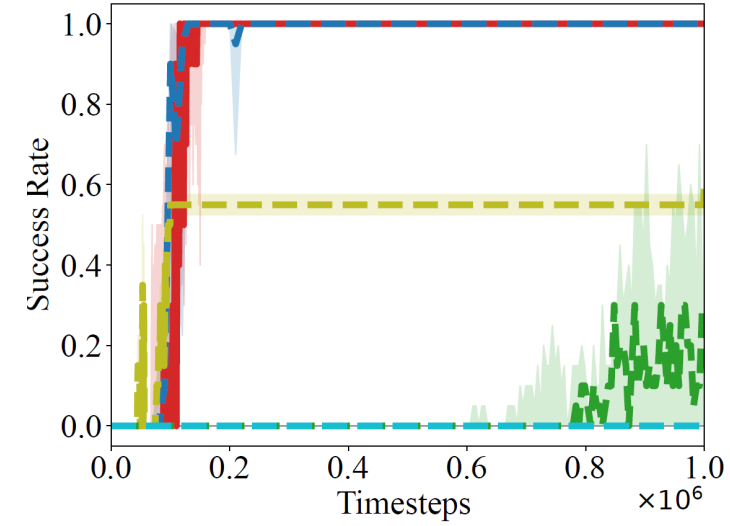


PointSpiralMaze

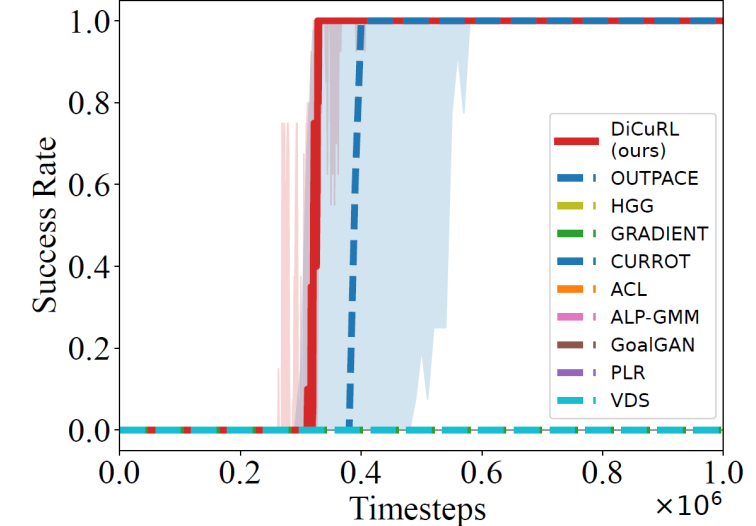
Results and Success Rate



PointUMaze



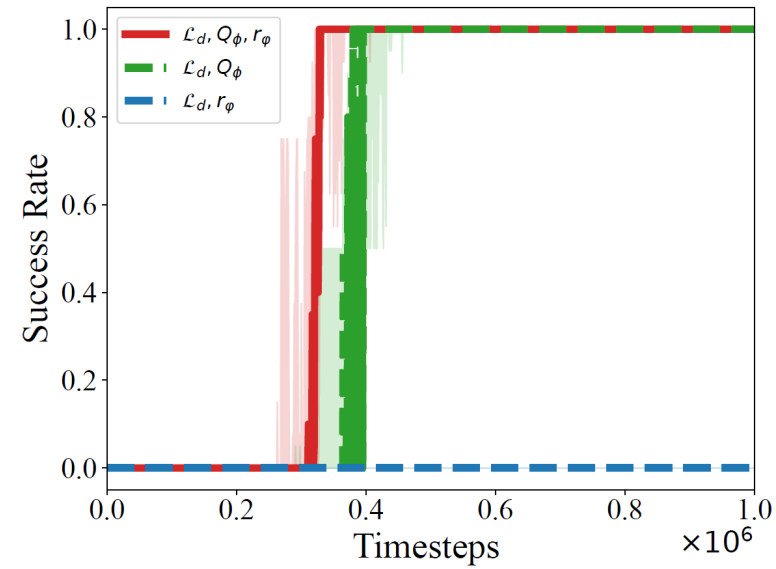
PointNMaze



PointSpiralMaze

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

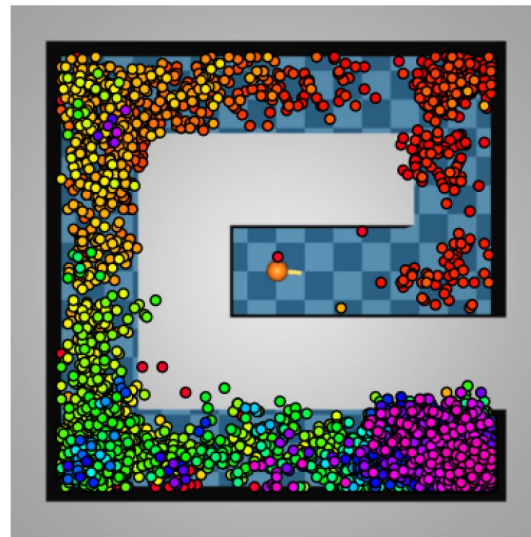
Ablation Study



$$\mathcal{L} = \xi_d \cdot \mathcal{L}_d(\theta) - \xi_q \cdot Q_\phi(s, g_0, \pi(s, g_0)) - \xi_r \cdot r_\phi(g_0, g_d)$$

$$\mathcal{L} = \xi_d \cdot \mathcal{L}_d(\theta) - \xi_q \cdot Q_\phi(s, g_0, \pi(s, g_0)) - \xi_r \cdot r_\phi(g_0, g_d)$$

$$\mathcal{L} = \xi_d \cdot \mathcal{L}_d(\theta) - \xi_q \cdot Q_\phi(s, g_0, \pi(s, g_0)) - \xi_r \cdot r_\phi(g_0, g_d)$$



Curriculum goals with Q



Curriculum goals with reward

Conclusions

- We introduced the **Diffusion Based Curriculum Reinforcement Learning** (DiCuRL) that utilizes diffusion models to generate curriculum goals.
- The diffusion model is trained to minimize its loss while simultaneously maximizing Q and Reward function.
- The generated goals promote exploration due to the inherent noising and denoising mechanism of the diffusion model.
- Additionally, we tested DiCuRL on two robot manipulation tasks, FetchPush and FetchPickAndPlace (see the results in the paper.)