
Enhancing the Hierarchical Environment Design via Generative Trajectory Modeling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Unsupervised Environment Design (UED) is a paradigm that automatically gener-
2 ates a curriculum of training environments, enabling agents trained in these
3 environments to develop general capabilities, i.e., achieving good zero-shot transfer
4 performance. However, existing UED approaches focus primarily on the random
5 generation of environments for open-ended agent training. This is impractical
6 in resource-limited scenarios where there is a constraint on the number of envi-
7 ronments that can be generated. In this paper, we introduce a hierarchical MDP
8 framework for environment design under resource constraints. It consists of an
9 upper-level RL teacher agent that generates suitable training environments for a
10 lower-level student agent. The RL teacher can leverage previously discovered
11 environment structures and generate environments at the frontier of the student’s
12 capabilities by observing the student policy’s representation. Additionally, to alle-
13 viate the time-consuming process of collecting the experience of the upper-level
14 teacher, we utilize recent advances in generative modeling to synthesize a trajec-
15 tory dataset for training the teacher agent. Our method significantly reduces the
16 resource-intensive interactions between agents and environments, and empirical
17 experiments across various domains demonstrate the effectiveness of our approach.

18 **1 Introduction**

19 The advances of reinforcement learning (RL) [17] have promoted research into the problem of
20 training autonomous agents that are capable of accomplishing complex tasks. One interesting, yet
21 underexplored, area is training agents to perform well in unseen environments, a concept referred to
22 as zero-shot transfer performance. To this end, Unsupervised Environment Design (UED) [3] has
23 emerged as a promising paradigm to address this problem. The objective of UED is to automatically
24 generate environments in a curriculum-based manner, and training agents in these sequentially
25 generated environments can equip agents with general capabilities, enabling agents to learn robust
26 and adaptive behaviors that can be transferred to new scenarios without explicit exposure during
27 training.

28 Existing approaches in UED primarily focus on building an adaptive curriculum for the environment
29 generation process to train the generally capable agent. Dennis et al. [3] formalize the problem of
30 finding adaptive curricula through a game involving an adversarial environment generator (teacher
31 agent), an antagonist agent (expert agent), and the protagonist agent (student agent). The RL-based
32 teacher is designed to generate environments that maximize regret, defined as the difference between
33 the protagonist and antagonist agent’s expected rewards. They show that these agents will reach
34 a Nash Equilibrium where the student agent learns the minimax regret policy. However, since the
35 teacher agent adapts solely based on the regret feedback, it is inherently difficult to adapt to student

36 policy changes. Meanwhile, training such an RL-based teacher remains a challenge because of the
37 high computational cost of training an expert antagonist agent for each environment.

38 In contrast, domain randomization [19] based approaches circumvent the overhead of developing
39 an RL teacher by training agents in randomly generated environments, resulting in good empirical
40 performances. Building upon this, Jiang et al. [7] introduce an emergent curriculum by sampling
41 randomly generated environments with high regret value ¹ to train the agent. Parker-Holder et al.
42 [10] then propose the adaptive curricula by manually designing a principled, regret-based curriculum,
43 which involves generating random environments with increasing complexity. While these domain
44 randomization-based algorithms have demonstrated good zero-shot transfer performance, they face
45 limitations in efficiently exploring large environment design spaces and exploiting the inherent
46 structure of previously discovered environments. Moreover, existing UED approaches typically
47 rely on open-ended learning, necessitating a long training horizon, which is unrealistic in the real
48 world due to resource constraints. Our goal is to develop a teacher policy capable of generating
49 environments that are perfectly matched to the current skill levels of student agents, thereby allowing
50 students to achieve optimal general capability within a strict budget for the number of environments
51 generated and within a shorter training time horizon.

52 In this paper, we address these challenges by introducing a novel, adaptive environment design
53 framework. The core idea involves using a hierarchical Markov Decision Process (MDP) to simul-
54 taneously formulate the evolution of an upper-level teacher agent, tasked with generating suitable
55 environments to train the lower-level student agent to achieve general capabilities. To accurately
56 guide the generation of environments at the frontier of the student agent’s current capabilities, we
57 propose approximating the student agent’s policy/capability by its performances across a set of diverse
58 evaluation environments, which acts as the state abstraction for the teacher’s decision-making process.
59 The transitions in the teacher’s state represent the trajectories of the student agent’s capability after
60 training in the generated environment. However, collecting experience for the upper-level teacher
61 agent is slow and resource-intensive, since each upper-level MDP transition evolves a complete
62 training cycle of the student agent on the generated environment. To accelerate the collection of
63 upper-level MDP experiences, we utilize advances in diffusion models that can generate new data
64 points capturing complex distribution properties, such as skewness and multi-modality, exhibited
65 in the collected dataset [11]. Specifically, we employ diffusion probabilistic model [15, 6] to learn
66 the evolution trajectory of student policy/capability and generate synthetic experiences to enhance
67 the training efficiency of the teacher agent. Our method, called *Synthetically-enhanced Hierarchical*
68 *Environment Design (SHED)*, automatically generates increasingly complex environments suited to
69 the current capabilities of student agents.

70 In summary, we make the following contributions:

- 71 • We develop a novel hierarchical MDP framework for UED that introduces a straightforward method
72 to represent the current capability level of the student agent.
- 73 • We introduce *SHED*, which utilizes diffusion-based techniques to generate synthetic experiences.
74 This method can accelerate the training of the off-policy teacher agent.
- 75 • We demonstrate that our method outperforms existing UED approaches (i.e., achieving a better
76 general capability under resource constraints) in different task domains.

77 2 Preliminaries

78 In this section, we provide an overview of two main research areas upon which our work is based.

79 2.1 Unsupervised Environment Design

The objective of UED is to generate a sequence of environments that effectively train the student agent to achieve a general capability. Dennis et al. [3] first model UED with an Underspecified Partially Observable Markov Decision Process (UPOMDP), which is a tuple

$$\mathcal{M} = \langle A, O, \Theta, S^{\mathcal{M}}, \mathcal{P}^{\mathcal{M}}, \mathcal{I}^{\mathcal{M}}, \mathcal{R}^{\mathcal{M}}, \gamma \rangle$$

¹They approximate the regret value by the Generalized Advantage Estimate [12].

80 . The UPOMDP has a set Θ representing the free parameters of the environments, which are
 81 determined by the teacher agent and can be distinct to generate the next new environment. Further,
 82 these parameters are incorporated into the environment-dependent transition function $\mathcal{P}^{\mathcal{M}} : S \times A \times$
 83 $\Theta \rightarrow S$. Here A represents the set of actions, S is the set of states. Similarly, $\mathcal{I}^{\mathcal{M}} : S \rightarrow O$ is the
 84 environment-dependent observation function, $\mathcal{R}^{\mathcal{M}}$ is the reward function, and γ is the discount factor.
 85 Specifically, given the environment parameters $\vec{\theta} \in \Theta$, we denote the corresponding environment
 86 instance as $\mathcal{M}_{\vec{\theta}}$. The student policy π is trained to maximize the cumulative rewards $V^{\mathcal{M}_{\vec{\theta}}}(\pi) =$
 87 $\sum_{t=0}^T \gamma^t r_t$ in the given environment $\mathcal{M}_{\vec{\theta}}$ under a time horizon T , and r_t are the collected rewards
 88 in $\mathcal{M}_{\vec{\theta}}$. Existing works on UED consist of two main strands: the RL-based environment generation
 89 approach and the domain randomization-based environment generation approach.

90 The RL-based generation approach was first formalized by Dennis et al. [3] as a self-supervised RL
 91 paradigm for generating environments. This approach involves co-evolving an environment generator
 92 policy (teacher) with an agent policy π (student), where the teacher’s role is to generate environment
 93 instances that best support the student agent’s continual learning. The teacher is trained to produce
 94 challenging yet solvable environments that maximize the regret measure, which is defined as the
 95 performance difference between the current student agent and a well-trained expert agent π^* within
 96 the current environment: $Regret^{\mathcal{M}_{\vec{\theta}}}(\pi, \pi^*) = V^{\mathcal{M}_{\vec{\theta}}}(\pi^*) - V^{\mathcal{M}_{\vec{\theta}}}(\pi)$.

97 The domain randomization-based generation approach, on the other hand, involves randomly generat-
 98 ing environments. Jiang et al. [7] propose to collect encountered environments with high learning
 99 potentials, which are approximated by the Generalized Advantage Estimation (GAE) [12], and then
 100 the student agent can selectively train in these environments, resulting in an emergent curriculum
 101 of increasing difficulty. Additionally, Parker-Holder et al. [10] adopt a different strategy by using
 102 predetermined starting points for the environment generation process and gradually increasing complex-
 103 ity. They manually divide the environment design space into different difficulty levels and employ
 104 human-defined edits to generate similar environments with high learning potentials. Their algorithm,
 105 ACCEL, is currently the state-of-the-art (SOTA) in the field, and we use an edited version of ACCEL
 106 as a baseline in our experiments.

107 2.2 Diffusion Probabilistic Models

108 Diffusion models [15] are a specific type of generative model that learns the data distribution.
 109 Recent advances in diffusion-based models, including Langevin dynamics and score-based generative
 110 models, have shown promising results in various applications, such as time series forecasting [18],
 111 robust learning [9], anomaly detection [21] as well as synthesizing high-quality images from text
 112 descriptions [8, 11]. These models can be trained using standard optimization techniques, such as
 113 stochastic gradient descent, making them highly scalable and easy to implement.

114 In a diffusion probabilistic model, we assume a d -dimensional random variable $x_0 \in \mathbb{R}^d$ with an
 115 unknown distribution $q(x_0)$. Diffusion Probabilistic model involves two Markov chains: a predefined
 116 forward chain $q(x_k|x_{k-1})$ that perturbs data to noise, and a trainable reverse chain $p_\phi(x_{k-1}|x_k)$ that
 117 converts noise back to data. The forward chain is typically designed to transform any data distribution
 118 into a simple prior distribution (e.g., standard Gaussian) by considering perturb data with Gaussian
 119 noise of zero mean and a fixed variance schedule $\{\beta_k\}_{k=1}^K$ for K steps:

$$q(x_k|x_{k-1}) = \mathcal{N}(x_k; \sqrt{1 - \beta_k}x_{k-1}, \beta_k \mathbf{I}) \quad \text{and} \quad q(x_{1:K}|x_0) = \prod_{k=1}^K q(x_k|x_{k-1}), \quad (1)$$

120 where $k \in \{1, \dots, K\}$, and $0 < \beta_{1:K} < 1$ denote the noise scale scheduling. As $K \rightarrow \infty$, x_K
 121 will converge to isotropic Gaussian noise: $x_K \rightarrow \mathcal{N}(0, \mathbf{I})$. According to the rule of the sum of
 122 normally distributed random variables, the choice of Gaussian noise provides a closed-form solution
 123 to generate arbitrary time-step x_k through:

$$x_k = \sqrt{\bar{\alpha}_k}x_0 + \sqrt{1 - \bar{\alpha}_k}\epsilon, \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}). \quad (2)$$

124 Here $\alpha_k = 1 - \beta_k$ and $\bar{\alpha}_k = \prod_{s=1}^k \alpha_s$. The reverse chain $p_\phi(x_{k-1}|x_k)$ reverses the forward process
 125 by learning transition kernels parameterized by deep neural networks. Specifically, considering the
 126 Markov chain parameterized by ϕ , denoising arbitrary Gaussian noise into clean data samples can be
 127 written as:

$$p_\phi(x_{k-1}|x_k) = \mathcal{N}(x_{k-1}; \mu_\phi(x_k, k), \Sigma_\phi(x_k, k)) \quad (3)$$

128 It uses the Gaussian form $p_\phi(x_{k-1}|x_k)$ because the reverse process has the identical function form as
 129 the forward process when β_t is small [15]. Ho et al. [6] consider the following parameterization of

130 $p_\phi(x_{k-1}|x_k)$:

$$\mu_\phi(x_k, k) = \frac{1}{\alpha_k} \left(x_k - \frac{\beta_k}{\sqrt{1-\alpha_k}} \epsilon_\phi(x_k, k) \right) \text{ and } \Sigma_\phi(x_k, k) = \tilde{\beta}_k^{1/2} \text{ where } \tilde{\beta}_k = \begin{cases} \frac{1-\alpha_{k-1}}{1-\alpha_k} \beta_k & k > 1 \\ \beta_1 & k = 1 \end{cases} \quad (4)$$

131 ϵ_ϕ is a trainable function to predict the noise vector ϵ from x_k . Ho et al. [6] show that training
 132 the reverse chain to maximize the log-likelihood $\int q(x_0) \log p_\phi(x_0) dx_0$ is equivalent to minimizing
 133 re-weighted evidence lower bound (ELBO) that fits the noise. They derive the final simplified
 134 optimization objective:

$$\mathcal{L}(\phi) = \mathbb{E}_{x_0, k, \epsilon} [\|\epsilon - \epsilon_\phi(\sqrt{\alpha_k}x_0 + \sqrt{1-\alpha_k}\epsilon, k)\|^2]. \quad (5)$$

135 Once the model is trained, new data points can be subsequently generated by first sampling a random
 136 vector from the prior distribution, followed by ancestral sampling through the reverse Markov chain
 137 in Equation 3.

138 3 Approach

139 In this section, we formally describe our method, Synthetically-enhanced *Hierarchical Environment*
 140 *Design (SHED)*, which is a novel framework for UED under resource constraints. The *SHED*
 141 incorporates two key components that differentiate it from existing UED approaches:

- 142 • A hierarchical MDP framework to generate suitable environments,
- 143 • A generative model to generate the synthetic trajectories.

144 *SHED* uses a hierarchical MDP framework where an RL teacher leverages the observed student’s
 145 policy representation to generate environments at the student’s capabilities frontier. Such targeted
 146 environment generation process enhances the student’s general capability by utilizing the underlying
 147 structure of previously discovered environments, rather than relying on the open-ended random
 148 generation. Besides, *SHED* leverages advances in generative models to generate synthetic trajectories
 149 that can be used to train the off-policy teacher agent, which significantly reduces the costly interactions
 150 between the agents and the environments. The overall framework is shown in Figure 1, and the
 151 pseudo-code is provided in Algorithm 1.

152 3.1 Hierarchical Environment Design

153 The objective is to generate a limited number of environments that are designed to enhance the general
 154 capability of the student agent. Inspired by the principles of PAIRED [3], we adopt an RL-based
 155 approach for the environment generation process. To better generate suitable environments tailored
 156 to the current student skill level, *SHED* uses the hierarchical MDP framework, consisting of an
 157 upper-level RL teacher policy Λ and a lower-level student policy π . Specifically, the teacher policy,
 158 $\Lambda : \Pi \rightarrow \Theta$, maps from the space of all potential student policies Π to the space of environment
 159 parameters Θ . Existing RL-based methods (e.g., PARIED) rely solely on regret feedback and
 160 fail to effectively capture the nuances of the student policy. To address this challenge, *SHED*
 161 enhances understanding by encoding the student policy π into a vector that serves as the state
 162 abstraction for teacher Λ . Rather than compressing the knowledge in the student policy network, we
 163 approximate the embedding of the student policy π by assessing performance across a set of diverse
 164 evaluation environments. This performance vector, denoted as $p(\pi)$, gives us a practical estimate
 165 of the student’s current general capabilities, enabling the teacher to customize the next training
 166 environments accordingly. In our hierarchical framework, the environment generation process is
 167 governed by discrete-time dynamics. We delve into the specifics below.

168 **Upper-level teacher MDP.** The upper-level teacher operates at a coarser layer of student policy
 169 abstraction and generates environments to train the lower-level student agent. This process can be
 170 formally modeled as an MDP by the tuple $\langle S^u, A^u, P^u, R^u, \gamma^u \rangle$:

- 171 • S^u represents the upper-level state space. Typically, $s^u = p(\pi) = [p_1, \dots, p_m]$ denotes the
 172 student performance vector across m diverse evaluation environments. This vector serves as the
 173 representation of the student policy π and is observed by the teacher.

Algorithm 1 *SHED*

Input: real data ratio $\psi \in [0, 1]$, evaluate environment set θ^{eval} , reward function R ;

- 1: **Initialize:** diffusion model D , teacher policy Λ , real and synthetic replay buffer $\mathcal{B}_{\text{real}}, \mathcal{B}_{\text{syn}} = \emptyset$;
- 2: **for** episode $ep = 1, \dots, K$ **do**
- 3: Initialize student policy π
- 4: Evaluate π on θ^{eval} and get state $s^u = p(\pi)$
- 5: **for** Budget $t = 1, \dots, T$ **do**
- 6: generate $\vec{\theta} \sim \Lambda$, and create $\mathcal{M}_{\vec{\theta}}(\pi)$
- 7: train π on $\mathcal{M}_{\vec{\theta}}$ to maximize $V^{\vec{\theta}}(\pi)$
- 8: evaluate π on θ^{eval} and get next state s'
- 9: compute teacher's reward r_t according to R
- 10: add experience $(s_t^u, \vec{\theta}, r_t, s_t^{u'})$ to $\mathcal{B}_{\text{real}}$
- 11: train D with samples from $\mathcal{B}_{\text{real}}$
- 12: generate synthetic experiences from D and add them to \mathcal{B}_{syn}
- 13: train Λ on samples from $\mathcal{B}_{\text{real}} \cup \mathcal{B}_{\text{syn}}$ mixed with ratio ψ
- 14: set $s = s'$;
- 15: **end for**
- 16: **end for**

Output: Λ, π, D

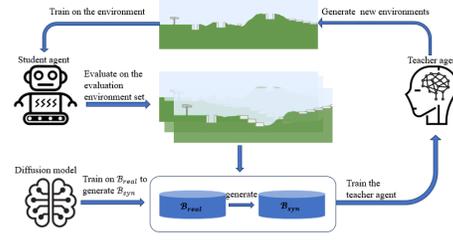


Figure 1: The overall framework of *SHED*.

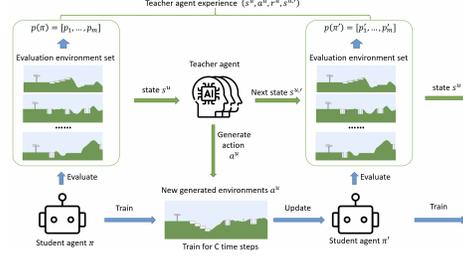


Figure 2: The illustration of the environment generation process.

174 • A^u is the upper-level action space. The teacher observes the abstraction of the student policy,
 175 s^u and produces an upper-level action a^u which is the environment parameters $\vec{\theta}$. $\vec{\theta}(a^u)$ is then
 176 used to generate specific environment instances $\mathcal{M}_{\vec{\theta}}$. Thus the upper-level action space A^u is the
 177 environment parameter space Θ .

178 • P^u denotes the action-dependent transition dynamics of the upper-level state. The general capability
 179 of the student policy evolves due to training the student agent on the generated environments.

180 • R^u provides the upper-level reward to the teacher at the end of training the student on the generated
 181 environment. The design of R^u will be discussed in Section 3.3.

182 As shown in Figure 2, given the student policy π , the teacher Λ first observes the representation
 183 of the student policy, $s^u = [p_1, \dots, p_m]$. Then teacher produces an upper-level action a^u which
 184 corresponds to the environment parameters. These environment parameters are subsequently used
 185 to generate specific environment instances. The lower-level student policy π will be trained on the
 186 generated environments for C training steps. The upper-level teacher collects and stores the student
 187 policy evolution transition $(s^u, a^u, r^u, s^{u'})$ every C times steps for off-policy training. The teacher
 188 agent is trained to maximize the cumulative reward giving the budget for the number of generated
 189 environments. The choice of the evaluation environments will be discussed in Section 3.3.

190 **Lower-level student MDP.** The generated environment is fully specified for the student, characterized
 191 by a Partially Observable Markov Decision Process (POMDP), which is defined by a tuple $\mathcal{M}_{\vec{\theta}} = \langle$
 192 $A, O, S^{\vec{\theta}}, \mathcal{P}^{\vec{\theta}}, \mathcal{I}^{\vec{\theta}}, \mathcal{R}^{\vec{\theta}}, \gamma \rangle$, where A represents the set of actions, O is the set of observations, $S^{\vec{\theta}}$
 193 is the set of states determined by the environment parameters $\vec{\theta}$, similarly, $\mathcal{P}^{\vec{\theta}}$ is the environment-
 194 dependent transition function, and $\mathcal{I}^{\vec{\theta}}: \vec{\theta} \rightarrow O$ is the environment-dependent observation function,
 195 $\mathcal{R}^{\vec{\theta}}$ is the reward function, and γ is the discount factor. At each time step t , the environment produces a
 196 state observation $s_t \in S^{\vec{\theta}}$, the student agent samples the action $a_t \sim A$ and interacts with environment
 197 $\vec{\theta}$. The environment yields a reward r_t according to the reward function $\mathcal{R}^{\vec{\theta}}$. The student agent is
 198 trained to maximize their cumulative reward $V^{\vec{\theta}}(\pi) = \sum_{t=0}^C \gamma^t r_t$ for the current environment under
 199 a finite time horizon C . The student agent will learn a good general capability from training on a
 200 sequence of generated environments.

201 The hierarchical framework enables the teacher agent to systematically measure and enhance the
 202 general capability of the student agent and to adapt the training process accordingly. However, it’s
 203 worth noting that collecting student policy evolution trajectories $(s^u, a^u, r^u, s^{u'})$ to train the teacher
 204 agent is notably slow and resource-intensive, since each transition in the upper-level teacher MDP
 205 encompasses a training horizon of C timesteps for the student in the generated environment. Thus, it
 206 is essential to reduce the need for costly collection of upper-level teacher experiences.

207 3.2 Generative Trajectory Modeling

208 In this section, we will formally introduce a generative model designed to ease the collection of upper-
 209 level MDP experience. This will allow us to train our teacher policy more efficiently. In particular, we
 210 first utilize a diffusion model to learn the conditional data distribution from the collected experiences
 211 $\tau = \{(s_t^u, a_t^u, r_t^u, s_t^{u'})\}$. Later we can use the reverse chain in the diffusion model to generate the
 212 synthetic trajectories that can be used to help train the teacher agent, thereby alleviating the need
 213 for extensive and time-consuming collection of upper-level teacher experiences. We deal with two
 214 different types of timesteps in this section: one for the diffusion process and the other for the upper-
 215 level teacher agent, respectively. We use subscripts $k \in 1, \dots, K$ to represent diffusion timesteps
 216 and subscripts $t \in 1, \dots, T$ to represent trajectory timesteps in the teacher’s experience.

217 In the image domain, the diffusion process is implemented across all pixel values of the image. In our
 218 setting, we diffuse over the next state $s^{u'}$ conditioned the given state s^u and action a^u . We construct
 219 our generative model according to the conditional diffusion process:

$$q(s_k^{u'} | s_{k-1}^{u'}), \quad p_\phi(s_{k-1}^{u'} | s_k^{u'}, s^u, a^u)$$

220 As usual, $q(s_k^{u'} | s_{k-1}^{u'})$ is the predefined forward noising process while $p_\phi(s_{k-1}^{u'} | s_k^{u'}, s^u, a^u)$ is the
 221 trainable reverse denoising process. We begin by randomly sampling the collected experiences
 222 $\tau = \{(s_t^u, a_t^u, r_t^u, s_t^{u'})\}$ from the real experience buffer \mathcal{B}_{real} . Giving the observed state s^u and
 223 action a^u , we use the reverse process p_ϕ to represent the generation of the next state $s^{u'}$:

$$p_\phi(s_{0:K}^{u'} | s^u, a^u) = \mathcal{N}(s_K^{u'}; 0, \mathbf{I}) \prod_{k=1}^K p_\phi(s_{k-1}^{u'} | s_k^{u'}, s^u, a^u)$$

At the end of the reverse chain, the sample $s_0^{u'}$, is the generated next state $s^{u'}$. Similar to Ho et al.
 [6], we parameterize $p_\phi(s_{k-1}^{u'} | s_k^{u'}, s^u, a^u)$ as a noise prediction model with the covariance matrix
 fixed as $\Sigma_\phi(s_k^{u'}, s^u, a^u, k) = \beta_k \mathbf{I}$, and the mean is

$$\mu_\phi(s_k^{u'}, s^u, a^u, k) = \frac{1}{\sqrt{\alpha_k}} \left(s_k^{u'} - \frac{\beta_k}{\sqrt{1 - \alpha_k}} \epsilon_\phi(s_k^{u'}, s^u, a^u, k) \right)$$

224 $\epsilon_\phi(s_k^{u'}, s^u, a^u, k)$ is the trainable denoising function, which aims to estimate the noise ϵ in the noisy
 225 input $s_k^{u'}$ at step k .

226 **Training objective.** We employ a similar simplified objective to train the conditional ϵ - model:

$$\mathcal{L}(\phi) = \mathbb{E}_{(s^u, a^u, s^{u'}) \sim \tau, k \sim \mathcal{U}, \epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon - \epsilon_\phi(s_k^{u'}, s^u, a^u, k)\|^2] \quad (6)$$

227 Where $s_k^{u'} = \sqrt{\alpha_k} s^{u'} + \sqrt{1 - \alpha_k} \epsilon$. The intuition for the loss function $\mathcal{L}(\phi)$ is to predict the noise
 228 $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ at the denoising step k , and the diffusion model is essentially learning the student
 229 policy involution trajectories collected in the real experience buffer \mathcal{B}_{reals} . Note that the reverse
 230 process necessitates a substantial number of steps K [15]. Recent research by Xiao et al. [22] has
 231 demonstrated that enabling denoising with large steps can reduce the total number of denoising steps
 232 K . To expedite the relatively slow reverse sampling process (as it requires computing ϵ_ϕ networks
 233 K times), we use a small value of K . Similar to Wang et al. [20], while simultaneously setting
 234 $\beta_{\min} = 0.1$ and $\beta_{\max} = 10.0$, we define:

$$\beta_k = 1 - \exp \left(\beta_{\min} \times \frac{1}{K} - 0.5(\beta_{\max} - \beta_{\min}) \frac{2k - 1}{K^2} \right)$$

235 This noise schedule is derived from the variance-preserving Stochastic Differential Equation by Song
 236 et al. [16].

237 **Generate synthetic trajectories.** Once the diffusion model has been trained, it can be used to generate
 238 synthetic experience data by starting with a draw from the prior $s_K^{u,\prime} \sim \mathcal{N}(0, \mathbf{I})$ and successively
 239 generating denoised next state, conditioned on the given s^u and a^u through the reverse chain p_ϕ .
 240 Note that the giving condition action a can either be randomly sampled from the action space or use
 241 another diffusion model to learn the action distribution giving the initial state s^u . This new diffusion
 242 model is essentially a behavior-cloning model that aims to learn the teacher policy $\Lambda(a^u|s^u)$. This
 243 process is similar to the work of Wang et al. [20]. We discuss this process in detail in the appendix.
 244 In this paper, we randomly sample a^u as it is straightforward and can also increase the diversity in
 245 the generated synthetic experience to help train a more robust teacher agent.

246 After obtaining the generated next state $s^{u,\prime}$ conditioned on s^u, a^u , we compute reward r^u using
 247 teacher’s reward function $R(s^u, a^u, s^{u,\prime})$. The specifics of how the reward function is chosen are
 248 explained in the following section.

249 3.3 Rewards and Choice of evaluate environments

250 **Selection of evaluation environments.** The upper-level teacher generates environments tailored
 251 for the lower-level student to improve its general capability. Thus it is important to select a set of
 252 diverse suitable evaluation environments as the performance vector reflects the student agent’s general
 253 capabilities and serves as an approximation of the policy’s embedding. Fontaine and Nikolaidis
 254 [5] propose the use of quality diversity (QD) optimization to collect high-quality environments that
 255 exhibit diversity for the agent behaviors. Similarly, Bhatt et al. [1] introduce a QD-based algorithm for
 256 dynamically designing such evaluation environments based on the current agent’s behavior. However,
 257 it’s worth noting that this QD-based approach can be tedious and time-consuming, and the collected
 258 evaluation environments heavily rely on the given agent policy.

259 Given these considerations, it is natural to take advantage of the domain randomization algorithm,
 260 as it has demonstrated compelling results in generating diverse environments and training generally
 261 capable agents. In our approach, we first discretize the environment parameters into different ranges,
 262 then randomly sample from these ranges, and combine these parameters to generate evaluation
 263 environments. This method can generate environments that may induce a diverse performance for the
 264 same policy, and it shows promising empirical results in the final experiments.

265 **Reward design.** We define the reward function for the upper-level teacher policy as a parameterized
 266 function based on the improvement in student performance in the evaluation environments after
 267 training in the generated environment:

$$R(s^u, a^u, s^{u,\prime}) = \sum_{i=1}^m (p'_i - p_i)$$

268 This reward function gives positive rewards to the upper-level teacher for taking action to create
 269 the right environment to improve the overall performance of students across diverse environments.
 270 However, it may encourage the teacher to obtain higher rewards by sacrificing student performance
 271 in one subset of evaluation environments to improve student performance in another subset, which
 272 conflicts with our objective to develop a student agent with general capabilities. Therefore, we need
 273 to consider fairness in the reward function to ensure that the generated environment can improve
 274 student’s general capabilities. Similar to [4], we build our fairness metric on top of the change
 275 in student’s performance in each evaluation environment, denoted as $\omega_i = p'_i - p_i$, and we have
 276 $\bar{\omega} = \frac{1}{m} \sum_{i=1}^m \omega_i$. We then measure the fairness of the teacher’s action using the coefficient of
 277 variation of student performances:

$$cv(s^u, a^u, s^{u,\prime}) = \sqrt{\frac{1}{m-1} \sum_i \frac{(\omega_i - \bar{\omega})^2}{\bar{\omega}^2}} \quad (7)$$

278 A teacher is considered to be fair if and only if the cv is smaller. As a result, our reward function is:

$$R(s^u, a^u, s^{u,\prime}) = \sum_{i=1}^m (p'_i - p_i) - \eta \cdot cv(s^u, a^u, s^{u,\prime}) \quad (8)$$

279 Here η is the coefficient that balances the weight of fairness in the reward function (We set a small
 280 value to η). This reward function motivates the teacher to generate training environments that can
 281 improve student’s general capability.

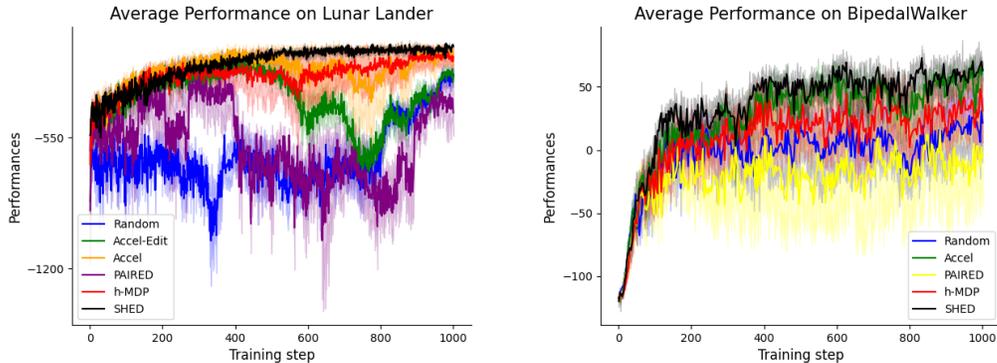


Figure 3: *Left*: The average zero-shot transfer performances on the test environments in the Lunar lander environment (mean and standard error). *Right*: The average zero-shot transfer performances on the test environments in the BipedalWalker (mean and standard error).

282 4 Experiments

283 In this section, we conduct experiments to compare *SHED* to other leading approaches on three
 284 domains: Lunar Lander, maze and a modified BipedalWalker environment. Experimental details and
 285 hyperparameters can be found in the Appendix. Specifically, our primary comparisons involve *SHED*
 286 and *h-MDP* (our proposed hierarchical approach without diffusion model aiding in training) against
 287 four baselines: domain randomization [19], ACCEL, [10], Edited ACCEL (with slight modifications
 288 that it does not revisit the previously generated environments), PAIRED [3]. In all cases, we
 289 train a student agent via Proximal Policy Optimization (PPO [13]), and train the teacher agent via
 290 Deterministic policy gradient algorithms (DDPG [14]), because DDPG is an off-policy algorithm and
 291 can learn from both real experiences and the synthetic experiences.

292 **Setup.** For each domain, we construct a set of evaluation environments and a set of test environments.
 293 The vector of student performances in the evaluation environments is used as the approximation of
 294 the student policy (as the observation to teacher agent), and the performances in the test environments
 295 are used to represent the student’s zero-shot transfer performances (general capabilities). Note that in
 296 order to obtain a fair comparison of zero-shot transfer performance, the evaluation environments and
 297 test environments do not share the same environment and they are not present during training.

298 **Lunar Lander.** This is a classic rocket trajectory optimization problem. In this domain, student
 299 agents are tasked with controlling a lander’s engine to safely land the vehicle. Before the start of each
 300 episode, teacher algorithms determine the environment parameters that are used to generate environ-
 301 ments in a given play-through, which includes gravity, wind power, and turbulence power. These
 302 parameters directly alter the difficulty of landing the vehicle safely. The state is an 8-dimensional
 303 vector, which includes the coordinates of the lander, its linear velocities, its angle, its angular velocity,
 304 and two booleans that represent whether each leg is in contact with the ground or not.

305 We train the student agent for $1e6$ environment time steps and periodically test the agent in test
 306 environments. The parameters for the test environments are randomly generated and fixed during
 307 training. We report the experiment results on the left side of Figure 3. As we can see, student
 308 agents trained under *SHED* consistently outperform other baselines and have minimal variance in
 309 transfer performance. During training, the baselines, except *h-MDP*, show a performance dip in the
 310 middle. This phenomenon could potentially be attributed to the inherent challenge of designing the
 311 appropriate environment instance in the large environment parameter space. This further demonstrates
 312 the effectiveness of our hierarchical design (*SHED* and *h-MDP*), which can successfully create
 313 environments that are appropriate to the current skill level of the students.

314 **Bipedalwalker.** We also evaluate *SHED* in the modified BipedalWalker from Parker-Holder et al.
 315 [10]. In this domain, the student agent is required to control a bipedal vehicle and navigate across the
 316 terrain, and the student receives a 24-dimensional proprioceptive state with respect to its lidar sensors,
 317 angles, and contacts. The teacher is tasked to select eight variables (including ground roughness, the

318 number of stairs steps, min/max range of pit gap width, min/max range of stump height, and min/max
319 range of stair height) to generate the corresponding terrain.

320 We use similar experiment settings in prior UED works, we train all the algorithms for $1e7$ environ-
321 ment time steps, and then evaluate their generalization ability on ten distinct test environments in
322 Bipedal-Walker domain. The parameters for the test environments are randomly generated and fixed
323 during training. As shown in Figure 3, our proposed method *SHED* surpasses all other baselines and
324 achieves performance levels nearly on par with the SOTA (ACCEL). Meanwhile, *SHED* maintains a
325 slight edge in terms of stability and overall performance and *PAIRED* suffers from a considerable
326 degree of variance in its performance.

327 **Partially observable Maze.** Here we study navigation tasks, where an agent must explore to find a
328 goal while navigating around obstacles. The environment is partially observable, and the agent’s field
329 of view is limited to a 3×3 grid area. Unlike the previously mentioned domains, maze environments
330 are non-parametric and cannot be directly represented by compact parameter vectors due to their
331 high complexity. To solve this challenge, we propose a novel method to generate maze by leveraging
332 advances in large language models (e.g., ChatGPT). Specifically, we implement a retrieval-augmented
333 generation (RAG) process to optimize the ChatGPT’s output such that it can generate desired maze
334 environments. This process ensures that large language models reference authoritative knowledge
335 bases to generate feasible mazes. To simplify the teacher’s action space, we extracted several key
336 factors that constitute the teacher’s action space (environmental parameters) for maze generation.
337 Details on maze generation are provided in Appendix D.3, and prompt are included in Appendix D.4.

338 The average zero-shot transfer perfor-
339 mances are reported in Figure 4. No-
340 tably, *SHED* demonstrates the highest
341 performance, consistently improving
342 and achieving the highest cumulative
343 rewards. The performance of h-MDP
344 steadily improves but does not reach
345 the highest levels, which further high-
346 lights the advantages of incorporat-
347 ing the generated synthetic datasets
348 to train an effective RL teacher agent.
349 Meanwhile, Accel-Edit and Accel
350 show higher variances in performance,
351 indicating that random teachers are
352 less stable in finding a suitable envi-
353 ronment to train student agents.

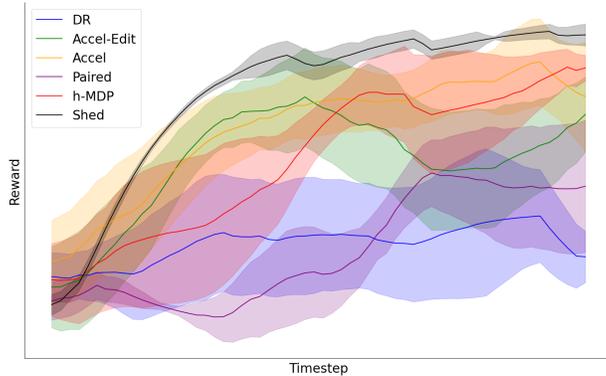


Figure 4: Average zero-shot transfer performance on the test environments in the maze environments.

354 **Ablation and additional Experi-**
355 **ments** In Appendix C, we evaluate
356 the ability of the diffusion model to generate the synthetic student policy involution trajectories. We
357 further provide ablation studies to assess the impact of different design choices in Appendix E.1.
358 Additionally, in Appendix E.2, we conduct experiments to show how the algorithm performs under
359 different settings, including scenarios with a larger budget constraint on the number of generated
360 environments or a larger weight assigned to CV fairness rewards. Notably, all results consistently
361 demonstrate the effectiveness of our approach.

362 5 Conclusion

363 In this paper, we introduce an adaptive approach for efficiently training a generally capable agent
364 under resource constraints. Our approach is general, utilizing an upper-level MDP teacher agent
365 that can guide the training of the lower-level MDP student agent agent. The hierarchical framework
366 can incorporate techniques from existing UED works, such as prioritized level replay (revisiting
367 environments with high learning potential). Furthermore, we have described a method to assist the
368 experience collection for the teacher when it is trained in an off-policy manner. Our experiment
369 demonstrates that our method outperforms existing UED methods, highlighting its effectiveness as a
370 curriculum-based learning approach within the UED framework.

References

- 371
- 372 [1] Varun Bhatt, Bryon Tjanaka, Matthew Fontaine, and Stefanos Nikolaidis. Deep surrogate
373 assisted generation of environments. *Advances in Neural Information Processing Systems*, 35:
374 37762–37777, 2022.
- 375 [2] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,
376 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative
377 interactive environments. *arXiv preprint arXiv:2402.15391*, 2024.
- 378 [3] Michael Dennis, Natasha Jaques, Eugene Vinitzky, Alexandre Bayen, Stuart Russell, Andrew
379 Critch, and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised
380 environment design. *Advances in neural information processing systems*, 33:13049–13061,
381 2020.
- 382 [4] Salma Elmalaki. Fair-iot: Fairness-aware human-in-the-loop reinforcement learning for har-
383 nassing human variability in personalized iot. In *Proceedings of the International Conference*
384 *on Internet-of-Things Design and Implementation*, pages 119–132, 2021.
- 385 [5] Matthew Fontaine and Stefanos Nikolaidis. Differentiable quality diversity. *Advances in Neural*
386 *Information Processing Systems*, 34:10040–10052, 2021.
- 387 [6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*
388 *in neural information processing systems*, 33:6840–6851, 2020.
- 389 [7] Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In *Interna-*
390 *tional Conference on Machine Learning*, pages 4940–4950. PMLR, 2021.
- 391 [8] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
392 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing
393 with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 394 [9] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar.
395 Diffusion models for adversarial purification. *arXiv preprint arXiv:2205.07460*, 2022.
- 396 [10] Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward
397 Grefenstette, and Tim Rocktäschel. Evolving curricula with regret-based environment design.
398 *arXiv preprint arXiv:2203.01302*, 2022.
- 399 [11] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton,
400 Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.
401 Photorealistic text-to-image diffusion models with deep language understanding. *Advances in*
402 *Neural Information Processing Systems*, 35:36479–36494, 2022.
- 403 [12] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-
404 dimensional continuous control using generalized advantage estimation. *arXiv preprint*
405 *arXiv:1506.02438*, 2015.
- 406 [13] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
407 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 408 [14] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
409 Deterministic policy gradient algorithms. In *International conference on machine learning*,
410 pages 387–395. Pmlr, 2014.
- 411 [15] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-
412 vised learning using nonequilibrium thermodynamics. In *International conference on machine*
413 *learning*, pages 2256–2265. PMLR, 2015.
- 414 [16] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
415 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*
416 *preprint arXiv:2011.13456*, 2020.
- 417 [17] Richard S Sutton, Andrew G Barto, et al. *Introduction to reinforcement learning*, volume 135.
418 MIT press Cambridge, 1998.

- 419 [18] Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. CSDI: Conditional score-based
420 diffusion models for probabilistic time series imputation. *Advances in Neural Information*
421 *Processing Systems*, 34:24804–24816, 2021.
- 422 [19] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel.
423 Domain randomization for transferring deep neural networks from simulation to the real world.
424 In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages
425 23–30. IEEE, 2017.
- 426 [20] Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive
427 policy class for offline reinforcement learning. In *The Eleventh International Conference on*
428 *Learning Representations*, 2023. URL <https://openreview.net/forum?id=AHvFDPi-FA>.
- 429 [21] Julian Wyatt, Adam Leach, Sebastian M Schmon, and Chris G Willcocks. Anoddpm: Anomaly
430 detection with denoising diffusion probabilistic models using simplex noise. In *Proceedings of*
431 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 650–656, 2022.
- 432 [22] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma
433 with denoising diffusion gans. *arXiv preprint arXiv:2112.07804*, 2021.

434 **A Theorem**

435 **Theorem 1** *There exists a finite evaluation environment set that can capture the student's general*
 436 *capabilities and the performance vector $[p_1, \dots, p_m]$ is a good representation of the student policy.*

437 To prove this, we first provide the following Assumption:

438 **Assumption 1** *Let $p(\pi, \vec{\theta})$ denote the performance of student policy π in an environment $\vec{\theta}$. For $\forall i$ -th*
 439 *dimension of the environment parameters, denoted as θ_i , when changing the θ_i to θ'_i to get a new*
 440 *environment $\vec{\theta}'$ while keeping other environment parameters fixed, there $\exists \delta_i > 0$, if $|\theta'_i - \theta_i| \leq \delta_i$, we*
 441 *have $|p(\pi, \vec{\theta}') - p(\pi, \vec{\theta})| \leq \epsilon_i$, where $\epsilon_i \rightarrow 0$.*

442 If this is true, we then can construct a finite set of environments, and the student performances in
 443 those environments can represent the performances in all potential environments generated within
 444 the certain environment parameters open interval combinations, and the set of those open intervals
 445 combinations cover the environment parameter space Θ .

446 We begin from the simplest case where we only consider using one environment parameter to generate
 447 environments, denoted as θ_i . We can construct a finite environment parameter set for environment
 448 parameters, which is $\{\theta_i^{min} + 1/2 * \delta_i, \theta_i^{min} + 3/2 * \delta_i, \theta_i^{min} + 5/2 * \delta_i, \dots, \theta_i^{max} - \delta_i/2\}$. Assume
 449 the set size is L_i . We let the set $\{\vec{\theta}_i\}_{i=1}^{L_i}$ denote the corresponding generated environments. This is
 450 served as the **representative environment set**. Then the student performances in those environments
 451 are denoted as $\{p(\pi, \vec{\theta}_i)\}_{i=1}^{L_i}$, which we call it as **representative performance vector set**. We can
 452 divide the space for θ_i into a finite set of open intervals with size L_i , which is $\{[\theta_i^{min}, \theta_i^{min} + 3/2 * \delta_i],$
 453 $(\theta_i^{min} + 1/2 * \delta_i, \theta_i^{min} + 5/2 * \delta_i), (\theta_i^{min} + 3/2 * \delta_i, \theta_i^{min} + 5/2 * \delta_i), \dots, (\theta_i^{max} - 3/2 * \delta_i, \theta_i^{max}]\}$,
 454 which we call it as **representative parameter interval set**, also denoted as $\{(\theta_i - \delta, \theta_i + \delta)\}_{i=1}^{L_i}$.
 455 For any environment generated in those intervals, denoted as $\vec{\theta}'_i$, the performance $p(\pi, \vec{\theta}'_i)$ can always
 456 be represented by the $p(\pi, \vec{\theta}_i)$ which is in the same interval, as $|p(\pi, \vec{\theta}'_i) - p(\pi, \vec{\theta}_i)| \leq \epsilon_i$, where
 457 $\epsilon_i \rightarrow 0$. In such cases, the finite set of environmental parameter intervals $\{\theta_i^{min} + 1/2 * \delta_i, \theta_i^{min} +$
 458 $3/2 * \delta_i, \theta_i^{min} + 5/2 * \delta_i, \dots, \theta_i^{max} - \delta_i/2\}$ fully covers the entire parameter space Θ . We can find
 459 a representative environment set $\{\vec{\theta}_i\}_{i=1}^{L_i}$ that is capable of approximating the performance of the
 460 student policy within the open parameter intervals combination. This set effectively characterizes the
 461 general performance capabilities of the student policy π .

462 Then we extend to two environment parameter design space cases. Let's assume that the environment
 463 is generated by two-dimension environment parameters. Then, for each environment parameter,
 464 $\theta_i \in \{\theta_{1i}, \theta_{2i}\}$. We can find the same open interval set for each parameter. Specifically, for each θ_i ,
 465 there exists a δ_i , such that if $|\theta'_i - \theta_i| \leq \delta_i$, we have $|p(\pi, \vec{\theta}') - p(\pi, \vec{\theta})| \leq \epsilon_i$, where $\epsilon_i \rightarrow 0$. Hence,
 466 we let $\delta = \min\{\delta_1, \delta_2\}$ and $\epsilon = \epsilon_1 + \epsilon_2$. Thus the new **representative environment set** is the set
 467 that includes the any combination of $\{\theta_{1i}, \theta_{2i}\}$ where $\theta_{1i} \in \{\vec{\theta}_j\}_{j=1}^{L_1}$ and $\theta_{2i} \in \{\vec{\theta}_j\}_{j=1}^{L_2}$. We can get
 468 the **representative performance vector set** as $\{p(\pi, [\vec{\theta}_i, \vec{\theta}_j])\}_{i \in [1, L_1], j \in [1, L_2]}$. We then can construct
 469 the **representative parameter interval set** as $\{[(\theta_i - \delta, \theta_i + \delta), (\theta_j - \delta, \theta_j + \delta)]\}_{i \in [1, L_1], j \in [1, L_2]}$.
 470 As a result, for any new environments $[\vec{\theta}'_i, \vec{\theta}'_j]$, we can find the representative environment whose
 471 environment parameters are in the same parameter interval $[\vec{\theta}_i, \vec{\theta}_j]$, such that their performance
 472 difference is smaller than $\epsilon = \epsilon_1 + \epsilon_2$ for all $\forall i \in [1, L_1], \forall j \in [1, L_2]$:

$$\begin{aligned}
 |p(\pi, [\vec{\theta}'_i, \vec{\theta}'_j]) - p(\pi, [\vec{\theta}_i, \vec{\theta}_j])| &= |p(\pi, [\vec{\theta}'_i, \vec{\theta}'_j]) - p(\pi, [\vec{\theta}_i, \vec{\theta}_j]) + p(\pi, [\vec{\theta}'_i, \vec{\theta}_j]) - p(\pi, [\vec{\theta}_i, \vec{\theta}_j])| \\
 &\leq |p(\pi, [\vec{\theta}'_i, \vec{\theta}'_j]) - p(\pi, [\vec{\theta}_i, \vec{\theta}_j])| + |p(\pi, [\vec{\theta}'_i, \vec{\theta}_j]) - p(\pi, [\vec{\theta}_i, \vec{\theta}_j])| \\
 &\leq \delta_j + \delta_i \\
 &= \delta
 \end{aligned}
 \tag{9}$$

473 In such cases, the finite set of environmental parameter intervals $\{[(\theta_i - \delta, \theta_i + \delta), (\theta_j - \delta, \theta_j +$
 474 $\delta)]\}_{i \in [1, L_1], j \in [1, L_2]}$ fully covers the entire parameter space Θ . We can find a representative environ-
 475 ment set $\{\vec{\theta}_i\}_{i=1}^{L_i}$ that is capable of approximating the performance of the student policy within the

Table 1: The teacher policies corresponding to the three approaches for UED. $U(\Theta)$ is a uniform distribution over environment parameter space, \tilde{D}_π is a baseline distribution, $\tilde{\theta}_\pi$ is the trajectory which maximizes regret of π , and v_π is the value above the baseline distribution that π achieves on that trajectory, c_π is the negative of the worst-case regret of π . Details are described in PAIRED [3].

UED Approaches	Teacher Policy	Decision Rule
DR [19]	$\Lambda(\pi) = U(\Theta)$	Randomly sample
PAIRED [3]	$\Lambda(\pi) = \{\tilde{\theta}_\pi : \frac{c_\pi}{v_\pi}, \tilde{D}_\pi : \text{otherwise}\}$	Minimax Regret
SHED (ours)	$\Lambda(\pi) = \arg \max_{\tilde{\theta} \in \Theta} Q_\pi(s = \pi, a = \tilde{\theta})$	Maximize reward

476 open parameter intervals combination. This set effectively characterizes the general performance
477 capabilities of the student policy π .

478 Similarly, we can show this still holds when the environment is constructed by a larger dimension
479 environment parameters, where we set $\delta = \min\{\delta_i\}$, and $\epsilon = \sum_i \epsilon_i$, and we have $\delta > 0, \epsilon \rightarrow 0$. The
480 overall logic is that we can find a finite set, which is called **representative environment set**, and
481 we can use performances in this set to represent any performances in the environments generated
482 in the **representative parameter interval set**, which is called **representative performance vector**
483 **set**. Finally, we can show that **representative parameter interval set** fully covers the environment
484 parameter space. Thus there exists a finite evaluation environment set that can capture the student’s
485 general capabilities and the performance vector, called **representative performance vector set**,
486 $[p_1, \dots, p_m]$ is a good representation of the student policy.

487 B Details about the Generative model

488 B.1 Generative model to generate synthetic next state

489 Here, we describe how to leverage the diffusion model to learn the conditional data distribution in the
490 collected experiences $\tau = \{(s_t^u, a_t^u, r_t^u, s_t^{u'})\}$. Later we can use the trainable reverse chain in the
491 diffusion model to generate the synthetic trajectories that can be used to help train the teacher agent,
492 resulting in reducing the resource-intensive and time-consuming collection of upper-level teacher
493 experiences. We deal with two different types of timesteps in this section: one for the diffusion
494 process and the other for the upper-level teacher agent, respectively. We use subscripts $k \in 1, \dots, K$
495 to represent diffusion timesteps and subscripts $t \in 1, \dots, T$ to represent trajectory timesteps in the
496 teacher’s experience.

497 In the image domain, the diffusion process is implemented across all pixel values of the image. In our
498 setting, we diffuse over the next state $s^{u'}$ conditioned the given state s^u and action a^u . We construct
499 our generative model according to the conditional diffusion process:

$$q(s_k^{u'} | s_{k-1}^{u'}), \quad p_\phi(s_{k-1}^{u'} | s_k^{u'}, s^u, a^u)$$

500 As usual, $q(s_k^{u'} | s_{k-1}^{u'})$ is the predefined forward noising process while $p_\phi(s_{k-1}^{u'} | s_k^{u'}, s^u, a^u)$ is the
501 trainable reverse denoising process. We begin by randomly sampling the collected experiences
502 $\tau = \{(s_t^u, a_t^u, r_t^u, s_t^{u'})\}$ from the real experience buffer \mathcal{B}_{real} .

503 We drop the superscript u here for ease of explanation. Given the observed state s and action a , we
504 use the reverse process p_ϕ to represent the generation of the next state s' :

$$p_\phi(s'_{0:K} | s, a) = \mathcal{N}(s'_K; 0, \mathbf{I}) \prod_{k=1}^K p_\phi(s'_{k-1} | s'_k, s, a) \quad (10)$$

At the end of the reverse chain, the sample s'_0 , is the generated next state s' . As shown in Section 2.2, $p_\phi(s'_{k-1} | s'_k, s, a)$ could be modeled as a Gaussian distribution $\mathcal{N}(s'_{k-1}; \mu_\theta(s'_k, s, a, k), \Sigma_\theta(s'_k, s, a, k))$. Similar to Ho et al. [6], we parameterize $p_\phi(s'_{k-1} | s'_k, s, a)$ as a noise prediction model with the covariance matrix fixed as

$$\Sigma_\theta(s'_k, s, a, k) = \beta_i \mathbf{I}$$

and mean is

$$\mu_\theta(s'_i, s, a, k) = \frac{1}{\sqrt{\alpha_k}} \left(s'_k - \frac{\beta_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(s'_k, s, a, k) \right)$$

505 Where $\epsilon_\theta(s'_k, s, a, k)$ is the trainable denoising function, which aims to estimate the noise ϵ in the
 506 noisy input s'_k at step k . Specifically, giving the sampled experience (s, a, s') , we begin by sampling
 507 $s'_K \sim \mathcal{N}(0, \mathbf{I})$ and then proceed with the reverse diffusion chain $p_\phi(s'_{k-1} | s'_k, s, a)$ for $k = K, \dots, 1$.
 508 The detailed expression for s'_{k-1} is as follows:

$$\frac{s'_k}{\sqrt{\alpha_k}} - \frac{\beta_k}{\sqrt{\alpha_k(1 - \bar{\alpha}_k)}} \epsilon_\theta(s'_k, s, a, k) + \sqrt{\beta_k} \epsilon, \quad (11)$$

509 where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. Note that $\epsilon = 0$ when $k = 1$.

510 **Training objective.** We employ a similar simplified objective, as proposed by Ho et al. [6] to train
 511 the conditional ϵ - model through the following process:

$$\mathcal{L}(\theta) = \mathbb{E}_{(s,a,s') \sim \tau, k \sim \mathcal{U}, \epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon - \epsilon_\phi(s'_k, s, a, k)\|^2] \quad (12)$$

512 Where $s'_k = \sqrt{\bar{\alpha}_k} s' + \sqrt{1 - \bar{\alpha}_k} \epsilon$. \mathcal{U} represents a uniform distribution over the discrete set $\{1, \dots, K\}$.
 513 The intuition for the loss function $\mathcal{L}(\theta)$ tries to predict the noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ at the denoising step k ,
 514 and the diffusion model is essentially learning the student policy involution trajectories collected in
 515 the real experience buffer \mathcal{B}_{reals} . Note that the reverse process necessitates a substantial number of
 516 steps K , as the Gaussian assumption holds true primarily under the condition of the infinitesimally
 517 limit of small denoising steps [15]. Recent research by Xiao et al. [22] has demonstrated that enabling
 518 denoising with large steps can reduce the total number of denoising steps K . To expedite the relatively
 519 slow reverse sampling process outlined in Equation 3.2 (as it requires computing ϵ_ϕ networks K
 520 times), we use a small value of K , while simultaneously setting $\beta_{\min} = 0.1$ and $\beta_{\max} = 10.0$.
 521 Similar to Wang et al. [20], we define:

$$\begin{aligned} \beta_k &= 1 - \alpha_k \\ &= 1 - \exp \left(\beta_{\min} \times \frac{1}{K} - 0.5(\beta_{\max} - \beta_{\min}) \frac{2k - 1}{K^2} \right) \end{aligned}$$

522 This noise schedule is derived from the variance-preserving Stochastic Differential Equation by Song
 523 et al. [16].

524 **Generate synthetic trajectories.** Once the diffusion model has been trained, it can be used
 525 to generate synthetic experience data by starting with a draw from the prior $s'_K \sim \mathcal{N}(0, \mathbf{I})$ and
 526 successively generating denoised next state, conditioned on the given s and a through the reverse
 527 chain p_ϕ in Equation 3.2. Note that the giving condition action a can either be randomly sampled
 528 from the action space (which is also the environment parameter space) or use another diffusion model
 529 to learn the action distribution giving the initial state s . In such case, this new diffusion model is
 530 essentially a behavior-cloning model that aims to learn the teacher policy $\Lambda(a|s)$. This process is
 531 similar to the work of Wang et al. [20]. We discuss this process in detail in the appendix. In this paper,
 532 we randomly sample a as it is straightforward and can also increase the diversity in the generated
 533 synthetic experience to help train a more robust teacher agent.

534 B.2 Generative model to generate synthetic action

535 Once the diffusion model has been trained, it can be used to generate synthetic experience data
 536 by starting with a draw from the prior $s'_K \sim \mathcal{N}(0, \mathbf{I})$ and successively generating denoised next
 537 state, conditioned on the given s and a through the reverse chain p_ϕ in Equation 3.2. Note that the
 538 giving condition action a can either be randomly sampled from the action space (which is also the
 539 environment parameter space) or we can train another diffusion model to learn the action distribution
 540 giving the initial state s , and then use the trained new diffusion model to sample the action a giving
 541 the state s . This process is similar to the work of Wang et al. [20].

542 In particular, We construct another conditional diffusion model as:

$$q(a_k | a_{k-1}), \quad p_\phi(a_{k-1} | a_k, s)$$

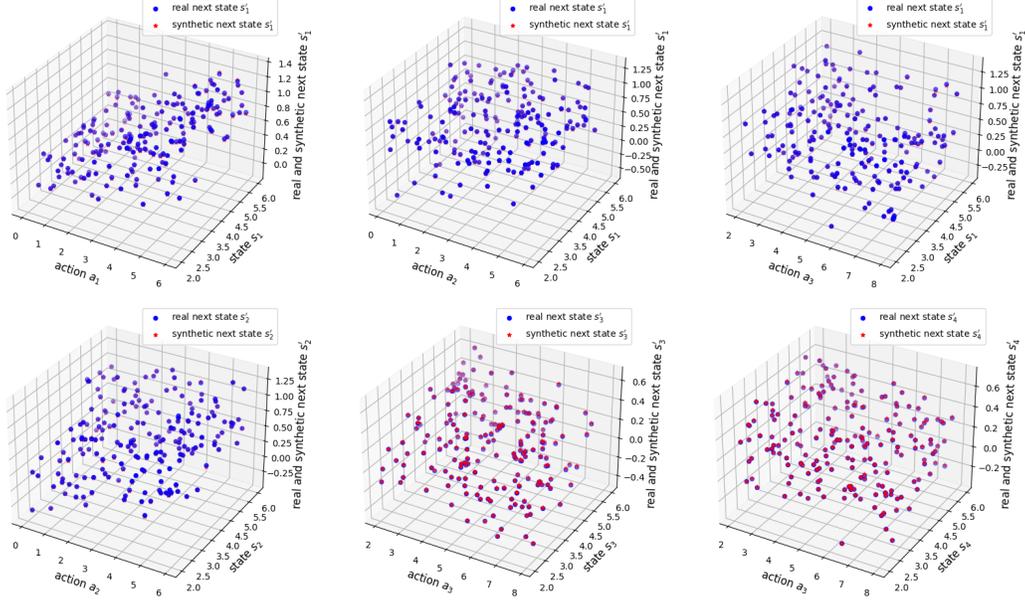


Figure 5: The distribution of the real s' and the synthetic s' conditioned on (s, a) .

543 As usual, $q(a_k|a_{k-1})$ is the predefined forward noising process while $p_\phi(a_{k-1}|a_k, s)$ is the trainable
 544 reverse denoising process. we represent the action generation process via the reverse chain of the
 545 conditional diffusion model as

$$p_\phi(a_{0:K}|s) = \mathcal{N}(a_K; 0, \mathbf{I}) \prod_{k=1}^K p_\phi(a_{k-1}|a_k, s) \quad (13)$$

At the end of the reverse chain, the sample a_0 , is the generated action a for the giving state s .
 Similarly, we parameterize $p_\phi(a_{k-1}|a_k, s)$ as a noise prediction model with the covariance matrix
 fixed as

$$\Sigma_\theta(a_k, s, k) = \beta_i \mathbf{I}$$

and mean is

$$\mu_\theta(a_i, s, k) = \frac{1}{\sqrt{\alpha_k}} \left(a_k - \frac{\beta_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(a_k, s, k) \right)$$

546 Similarly, the simplified loss function is

$$\mathcal{L}^a(\theta) = \mathbb{E}_{(s,a) \sim \tau, k \sim \mathcal{U}, \epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon - \epsilon_\phi(a_k, s, k)\|^2] \quad (14)$$

547 Where $a_k = \sqrt{\bar{\alpha}_k} a + \sqrt{1 - \bar{\alpha}_k} \epsilon$. \mathcal{U} represents a uniform distribution over the discrete set $\{1, \dots, K\}$.
 548 The intuition for the loss function $\mathcal{L}^a(\theta)$ tries to predict the noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ at the denoising step k ,
 549 and the diffusion model is essentially a behavior cloning model to learn the student policy collected
 550 in the real experience buffer \mathcal{B}_{reals} .

551 Once this new diffusion model is trained, the generation of the synthetic experience can be formulated
 552 as:

- 553 • we first randomly sample the state from the collected real trajectories $s \sim \tau$;
- 554 • we use the new diffusion model discussed above to mimic the teacher's policy to generate
 555 the actions a ;
- 556 • giving the state s and action a , we use the first diffusion model presented in the main paper
 557 to generate the next state s' ;
- 558 • we compute the reward r according to the reward function, and add the final generated
 559 synthetic experience (s, a, r, s') to the synthetic experience buffer \mathcal{B}_{syn} to help train the
 560 teacher agent.

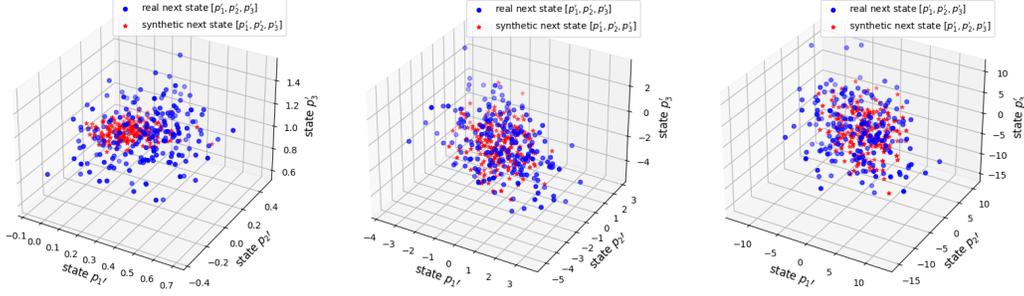


Figure 6: The distribution of the real $[s'_1, s'_2, s'_3]$ (red) and the synthetic $[s'_1, s'_2, s'_3]$ (blue) giving the fixed (s^u, a^u) . Specifically, the noise ε in $f(s^u, a^u)$ is (i).*left* figure: $\varepsilon = \epsilon$, (ii).*middle* figure: $\varepsilon = 3 * \epsilon$, (iii).*right* figure: $\varepsilon = 10 * \epsilon$, where $\epsilon \sim \mathcal{N}(0, 1)$.

561 C Empirical analysis of generative model

562 C.1 Ability to generate good synthetic trajectories

563 We begin by investigating *SHED*'s ability to assist in collecting experiences for the upper-level MDP
 564 teacher. This involves the necessity for *SHED* to prove its ability to accurately generate synthetic
 565 experiences for teacher agents. To check the quality of these generated synthetic experiences, we
 566 employ a diffusion model to simulate some data for validation (even though Diffusion models have
 567 demonstrated remarkable success across vision and NLP tasks).

568 We design the following experiment: given the teacher's observed state $s^u = [p_1, p_2, p_3, p_4, p_5]$,
 569 where p_i denotes the student performance on i -th evaluation environment. and given the teacher's
 570 action $a^u = [a_1, a_2, a_3]$, which is the environment parameters and are used to generate corresponding
 571 environment instances. We use a neural network $f(s^u, a^u)$ to mimic the involution trajectories of
 572 the student policy π . That is, with the input of the state s^u and action a^u into the neural network, it
 573 outputs the next observed state $s^{u,'} = [p'_1, p'_2, p'_3, p'_4, p'_5]$, indicating the updated student performance
 574 vector on the evaluation environments after training in the environment generated by a^u . In particular,
 575 we add a noise ε into $s^{u,'}$ to represent the uncertainty in the transition. We first train our diffusion
 576 model on the real dataset $(s^u, a^u, s^{u,'})$ generated by neural network $f(s^u, a^u)$. We then set a fixed
 577 (s^u, a^u) pair and input them into $f(s^u, a^u)$ to generate 200 samples of real $s^{u,'}$. The trained diffusion
 578 model is then used to generate 200 synthetic $s^{u,'}$ conditioned on the fixed (s^u, a^u) pair.

579 The results are presented in Figure 6, we can see that the generative model can effectively capture
 580 the distribution of real experience even if there is a large uncertainty in the transition, indicated by
 581 the value of ε . This provides evidence that the diffusion model can generate useful experiences
 582 conditioned on (s^u, a^u) . It is important to note that the marginal distribution derived from the reverse
 583 diffusion chain provides an implicit, expressive distribution, such distribution has the capability to
 584 capture complex distribution properties, including skewness and multi-modality.

585 C.2 addition experiments on diffusion model

586 We further provide more results to show the ability of our generative model to generate synthetic
 587 trajectories where the noise is extremely small. In such cases, the actual next state s' will converge to
 588 a certain value, and the synthetic next state $s^{sym,'}$ generated by the diffusion model should also be
 589 very close to that value, then the diffusion model has the ability to sample the next state $s^{sym,'}$ which
 590 can accurately represent the next state. We present the results in Figure 5. Specifically, this figure
 591 shows when the noise is very small in the actual next state, which is $0.05 * \epsilon$, and $\epsilon \sim \mathcal{N}(0, 1)$. Giving
 592 any condition (s, a) pair, we selectively report on (s_i, a_i) , where x -axis is the a_i value, and y -axis
 593 is the s_i value. The student policy with initial performance vector s is trained on the environments
 594 generated by the teacher's action a . We report the new performance s'_i of student policy on i -th
 595 environments after training in the z -axis. In particular, if two points s'_i and $s^{sym,'}_i$ are close, it indicates
 596 that the diffusion model can successfully generate the actual next state. As we can see, when the
 597 noise is extremely small, our diffusion model can accurately predict the next state of s'_i giving any
 598 condition (s, a) pair.

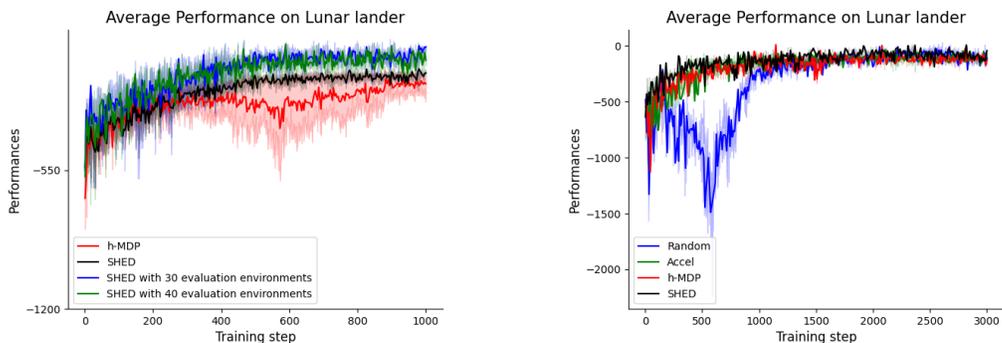


Figure 7: *Left*: The ablation study in the Lunar lander environment which investigates the effect of the size of the evaluation environment set. We provide the average zero-shot transfer performances on the test environments (mean and standard error). *Right*: Zero-shot transfer performance on the test environments under a longer time horizon in Lunar lander environments(mean and standard error).

599 D Additional Experiment Details

600 D.1 Hyperparameters

601 We set the learning rate $1e-3$ for actor, and $3e-3$ for critic, we set gamma $\gamma = 0.999$, $\lambda = 0.95$, and
 602 set coefficient for the entropy bonus (to encourage exploration) as 0.01. For each environment, we
 603 conduct 50 PPO updates for the student agent, and We can train on up to 50 environments, including
 604 replay. For our diffusion model, the diffusion discount is 0.99, and batch size is 64, τ is 0.005,
 605 learning rate is $3e-4$. The synthetic buffer size is 1000, and the ratio is 0.25.

606 D.2 Experiments Compute Resources

607 All the models were trained on a single NVIDIA GeForce RTX 3090 GPU and 16 CPUs.

608 D.3 Maze document

609 Here we provide the document shows the instruction to generate feasible maze environments.

610 There are several factors that can affect the difficulty of a maze. Here are
 611 some key factors to consider:

- 612 1. Maze Size: Larger mazes generally increase the complexity and difficulty
 613 as the agent has more states to explore. Typically, the maze size should be
 614 larger than 4x4 and smaller than 15*15.
 - 615 - If the size is 7*7 or smaller, the maze size is considered easy.
 - 616 - If the size is larger than 7*7 but smaller than 10*10, the maze size is
 617 considered medium.
 - 618 - If the maze size is larger than 10x10 but smaller than 15*15, the maze
 619 size is considered hard.
- 620 2. Maze Structure: The complexity of the paths, including the number of twists,
 621 turns, and dead-ends, can significantly impact navigation strategies. The
 622 presence of narrow corridors versus wide-open spaces also plays a role.
 - 623 - If there are fewer than 2 turns in the feasible path from the start position
 624 to the end position, the maze structure is considered easy.
 - 625 - If there are more than 2 turns but fewer than 4 turns in the path from the
 626 start position to the end position, the maze structure is considered medium.
 - 627 - If there are 4 or more turns in the path from the start position to the end
 628 position, the maze structure is considered hard.
- 629 3. Goal Location: The distance from the starting position to the end position
 630 also affects difficulty.
 - 631 - If the path from the start position to the end position requires fewer than

632 5 steps, the goal location is considered easy.
633 - If the path from the start position to the end position requires 5 to 10
634 steps, the goal location is considered medium.
635 - If the path from the start position to the end position requires more than
636 10 steps, the goal location is considered hard.
637 4. Start Location: The starting position can also affect the difficulty of
638 the maze. The starting position is categorized into five levels:
639 - If the start position is close to 1, it means it should be located as close
640 to the top left of the maze.
641 - If the start position is close to 2, it means it should be located as close
642 to the top right of the maze.
643 - If the start position is close to 3, it means it should be located as close
644 to the bottom left of the maze.
645 - If the start position is close to 4, it means it should be located as close
646 to the bottom right of the maze.
647 - If the start position is close to 5, it means it should be located as close
648 to the center of the maze.
649 Please note that the generated maze uses -1 to represent blocks, 0 to
650 represent the feasible path, 1 to represent the start position, and 2 to represent
651 the end position. Must ensure that there is a feasible path in the generated maze!
652 A feasible path means that 1 and 2 are connected directly through 0s, or 1 and 2
653 are connected directly. For example:
654 Feasible Maze:
655 Maze = [
656 [0, -1, -1, 2],
657 [1, -1, 0, 0],
658 [0, -1, 0, -1],
659 [0, 0, 0, -1],
660]
661 Non-Feasible Mazes:
662 Maze = [
663 [0, -1, -1, 2],
664 [1, -1, 0, 0],
665 [0, -1, -1, 0],
666 [0, 0, 0, -1],
667]
668 Or
669 Maze = [
670 [1, -1],
671 [-1, 2]
672]
673 These second example does not have any feasible path.
674
675

676 **D.4 Prompt for RAG**

677 We provide our prompt for the Retrieval Augmented Generation as follows:

678 Please refer to the document, and generate a maze with feasible path. The
679 difficulty level for the maze size is {maze_size_level}, and the difficulty
680 level for the maze structure is {maze_structure_level}, he difficulty level
681 for the goal location is {goal_location_level}, he difficulty level for
682 the start location is {start_position_level}.

683 E Additional experiments

684 E.1 Additional experiments about ablation studies

685 We also provide ablation analysis to evaluate the impact of different design choices in Lunar lander
686 domain, including (a) a larger evaluation environment set; (b) a bigger budget for constraint on the
687 number of generated environments (which incurs a longer training time horizon). The results are
688 reported in Figure 7.

689 We explore the impact of introducing the diffusion model in collecting synthetic teacher’s experience
690 and varying the size of the evaluation environment set. Specifically, as we can see from the right side
691 of Figure 7, the *SHED* consistently outperforms h-MDP, indicating the effectiveness of introducing
692 the generative model to help train the upper-level teacher policy. Furthermore, we find that when
693 increasing the size of the evaluation environment set, we can have a better result in the student
694 transfer performances. The intuition is that a larger evaluation environment set, encompassing a more
695 diverse range of environments, provides a better approximation of the student policy according to the
696 Theorem 1. However, the reason why *SHED* with 30 evaluation environments slightly outperforms
697 *SHED* with 40 evaluation environments is perhaps attributed to the increase in the dimension of the
698 student performance vector, which amplifies the challenge of training an effective diffusion model
699 with a limited dataset.

700 We conduct experiments in Lunar lander under a longer time horizon. The results are provided on the
701 right side of Figure 7. As we can see, our proposed algorithm *SHED* can efficiently train the student
702 agent to achieve the general capability in a shorter time horizon, This observation indicates that
703 our proposed environment generation process can better generate the suitable environments for the
704 current student policy, thereby enhancing its general capability, especially when there is a constraint
705 on the number of generated environments.

706 E.2 Additional experiments on Lunar lander

707 we also conduct experiments to show how the algorithm performs under different settings, such
708 as a larger weight of cv fairness rewards ($\eta = 10$). The results are provided in Figure 8. We
709 noticed an interesting finding: when fairness reward has a high weightage, our algorithm tends to
710 generate environments at the onset that lead to a rapid decline and subsequent improvement in student
711 performance across all test environments. This is done to avoid acquiring a substantial negative
712 fairness reward and thereby maximize the teacher’s cumulative reward. Notably, the student’s final
performance still surpasses other baselines at the end of training.

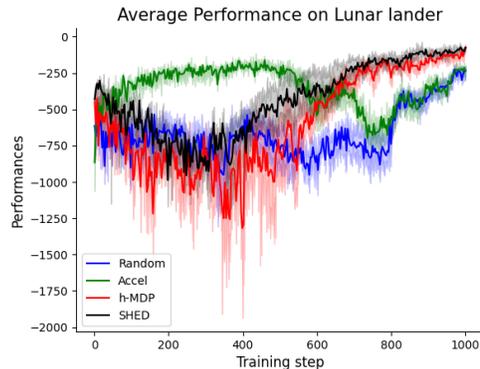


Figure 8: Zero-shot transfer performance on the test environments with a larger *cv* value coefficient in Lunar lander environments.

713

714 We further show in detail how the performance of different methods changes in each testing environ-
715 ment during training (see Figure 9 and Figure 10).

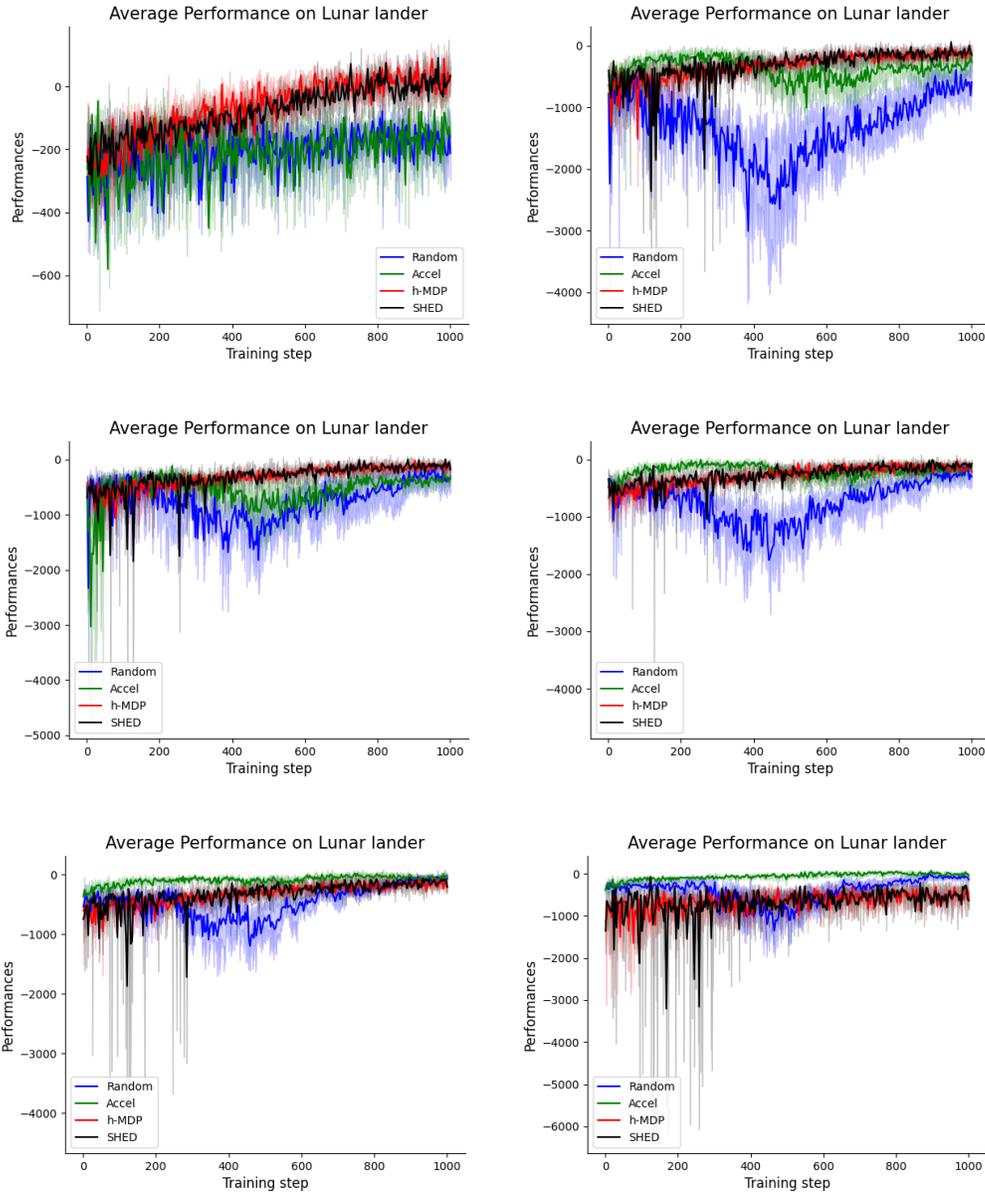


Figure 9: Detail how the performance of different methods changes in each testing environment during training (mean and error)

716 **E.3 Additional experiments on Maze**

717 We selectively report some results of zero-shot transfer performances in maze environments. The
 718 results are provided in Figure

719 **F Discussion**

720 **F.1 Limitations**

721 The limitation of this work comes from the UED framework, as UED is limited to the use of
 722 parameterized environments. This results in our experimental domain being relatively simple.

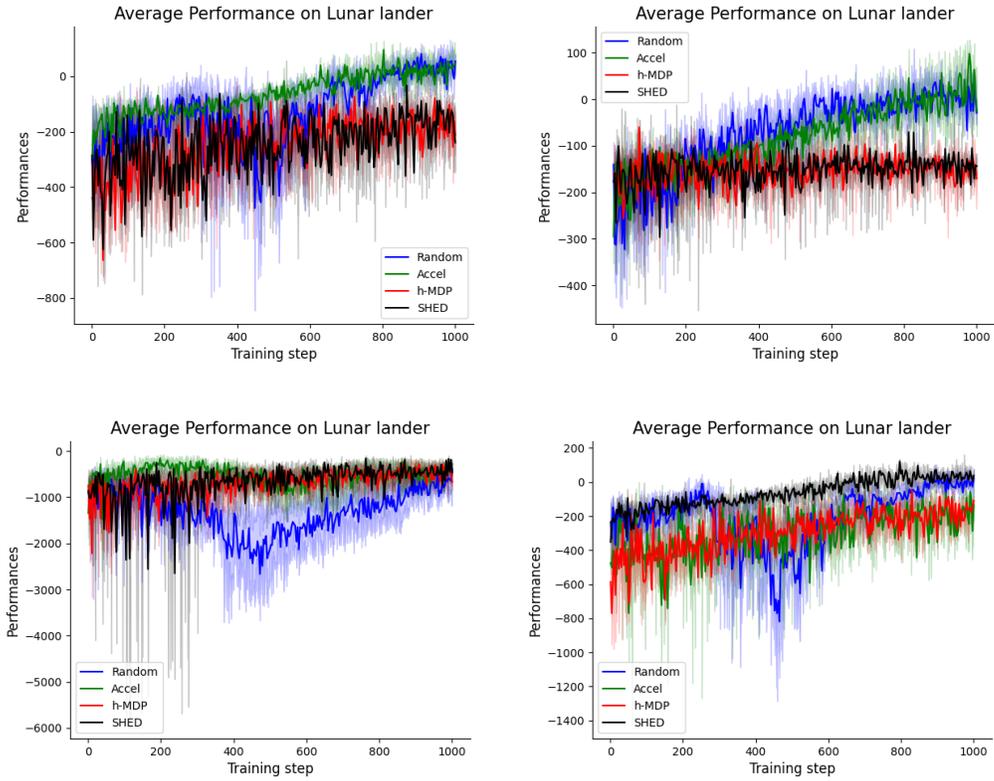


Figure 10: Detail how the performance of different methods changes in each testing environment during training (mean and error)

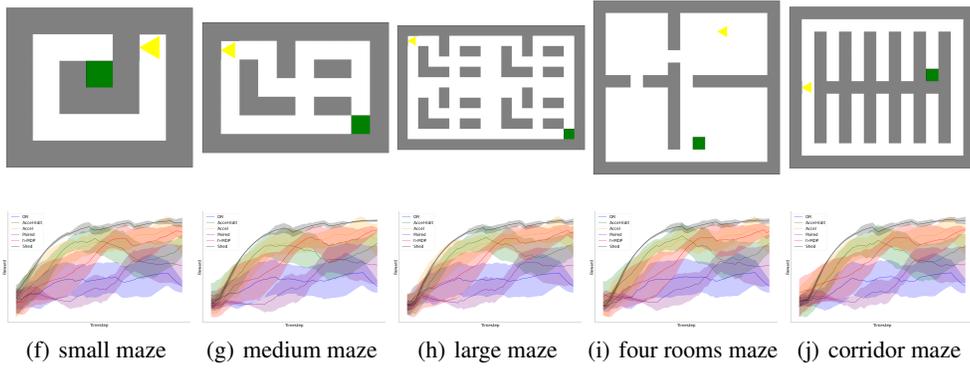


Figure 11: Zeros-shot transfer performance on test environments in maze environments

723 However, our work proposes a new hierarchical structure, and our policy representation is not only of
 724 great help for UED, but also has certain inspirations for hierarchical RL. Additionally, in the world
 725 model of UED (Genie [2]), the environment generator (teacher) focuses on creating video games, a
 726 domain that is compatible with our proposed application of upsampling the teacher agent’s experience
 727 using a diffusion model (since the state is image-based).

728 **NeurIPS Paper Checklist**

729 **1. Claims**

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

732 Answer: [Yes]

733 Justification: Yes, the main claims made in the abstract and introduction accurately reflect
734 the paper's contributions and scope.

735 Guidelines:

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

745 **2. Limitations**

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

747 Answer: [Yes]

748 Justification: The limitations of this work is discussed in Appendix F.1.

749 Guidelines:

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

776 **3. Theory Assumptions and Proofs**

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

779 Answer: [Yes]

780 Justification: See the theoretical result in Appendix 1.

781 Guidelines:

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

792 4. Experimental Result Reproducibility

793 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

794 perimental results of the paper to the extent that it affects the main claims and/or conclusions

795 of the paper (regardless of whether the code and data are provided or not)?

796 Answer: [Yes]

797 Justification: We disclose all the information needed to reproduce the main experimental

798 results of the paper to the extent that it affects the main claims and conclusions of the paper,

799 detailed in Section 3 and Appendix D.1.

800 Guidelines:

- 801 • The answer NA means that the paper does not include experiments.
- 802 • If the paper includes experiments, a No answer to this question will not be perceived
- 803 well by the reviewers: Making the paper reproducible is important, regardless of
- 804 whether the code and data are provided or not.
- 805 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 806 to make their results reproducible or verifiable.
- 807 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 808 For example, if the contribution is a novel architecture, describing the architecture fully
- 809 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 810 be necessary to either make it possible for others to replicate the model with the same
- 811 dataset, or provide access to the model. In general, releasing code and data is often
- 812 one good way to accomplish this, but reproducibility can also be provided via detailed
- 813 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 814 of a large language model), releasing of a model checkpoint, or other means that are
- 815 appropriate to the research performed.
- 816 • While NeurIPS does not require releasing code, the conference does require all submis-
- 817 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 818 nature of the contribution. For example
 - 819 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 820 to reproduce that algorithm.
 - 821 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 822 the architecture clearly and fully.
 - 823 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 824 either be a way to access this model for reproducing the results or a way to reproduce
 - 825 the model (e.g., with an open-source dataset or instructions for how to construct
 - 826 the dataset).
 - 827 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 828 authors are welcome to describe the particular way they provide for reproducibility.
 - 829 In the case of closed-source models, it may be that access to the model is limited in
 - 830 some way (e.g., to registered users), but it should be possible for other researchers
 - 831 to have some path to reproducing or verifying the results.

832 5. Open access to data and code

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

836 Answer: [Yes]

837 Justification: The code is provided in the supplementary marterial.

838 Guidelines:

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

858 6. Experimental Setting/Details

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

862 Answer: [Yes]

863 Justification: We provide the training and test details Section 3 and Appendix D.1 and
864 Appendix D.2.

865 Guidelines:

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

871 7. Experiment Statistical Significance

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

874 Answer: [Yes]

875 Justification: The proposed method is thoroughly evaluated on three domains, and the results
876 are reported based on a statistical analysis in Section 4 and Appendix E.

877 Guidelines:

- 878 • The answer NA means that the paper does not include experiments.
- 879 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
880 dence intervals, or statistical significance tests, at least for the experiments that support
881 the main claims of the paper.
- 882 • The factors of variability that the error bars are capturing should be clearly stated (for
883 example, train/test split, initialization, random drawing of some parameter, or overall
884 run with given experimental conditions).

- 885 • The method for calculating the error bars should be explained (closed form formula,
886 call to a library function, bootstrap, etc.)
- 887 • The assumptions made should be given (e.g., Normally distributed errors).
- 888 • It should be clear whether the error bar is the standard deviation or the standard error
889 of the mean.
- 890 • It is OK to report 1-sigma error bars, but one should state it. The authors should
891 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
892 of Normality of errors is not verified.
- 893 • For asymmetric distributions, the authors should be careful not to show in tables or
894 figures symmetric error bars that would yield results that are out of range (e.g. negative
895 error rates).
- 896 • If error bars are reported in tables or plots, The authors should explain in the text how
897 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

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

902 Answer: [Yes]

903 Justification: The detailed configuration of the experiments is listed with required computa-
904 tional resources.

905 Guidelines:

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

9. Code Of Ethics

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

917 Answer: [Yes]

918 Justification: We confirm that the research conducted in the paper conform, in every respect,
919 with the NeurIPS Code of Ethics, and all the authors preserve anonymity.

920 Guidelines:

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

10. Broader Impacts

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

929 Answer: [Yes]

930 Justification: The broader impacts of our paper are presented in Section F.1.

931 Guidelines:

- 932 • The answer NA means that there is no societal impact of the work performed.
- 933 • If the authors answer NA or No, they should explain why their work has no societal
934 impact or why the paper does not address societal impact.

- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

954 **11. Safeguards**

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

958 Answer: [NA]

959 Justification: Our paper poses no such risks.

960 Guidelines:

- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

971 **12. Licenses for existing assets**

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

975 Answer: [Yes]

976 Justification: All the assets, used in our paper, are properly credited and we explicitly
977 mention and properly respect the license and terms of use.

978 Guidelines:

- 979
- 980
- 981
- 982
- 983
- 984
- 985
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 986
- 987
- 988
- 989
- 990
- 991
- 992
- 993
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

994 **13. New Assets**

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

997 Answer: [NA]

998 Justification: This paper does not release new assets.

999 Guidelines:

- 1000
- 1001
- 1002
- 1003
- 1004
- 1005
- 1006
- 1007
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

1008 **14. Crowdsourcing and Research with Human Subjects**

1009 Question: For crowdsourcing experiments and research with human subjects, does the paper
1010 include the full text of instructions given to participants and screenshots, if applicable, as
1011 well as details about compensation (if any)?

1012 Answer: [NA]

1013 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1014 Guidelines:

- 1015
- 1016
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

1023 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
1024 Subjects**

1025 Question: Does the paper describe potential risks incurred by study participants, whether
1026 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1027 approvals (or an equivalent approval/review based on the requirements of your country or
1028 institution) were obtained?

1029 Answer: [NA]

1030 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1031 Guidelines:

- 1032
- 1033
- 1034
- 1035
- 1036
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

1037
1038
1039
1040
1041

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.