UNSUPERVISED REINFORCEMENT LEARNING BY MAXIMIZING SKILL DENSITY DEVIATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Unsupervised Reinforcement Learning (RL) aims to discover diverse behaviors that can accelerate the learning of downstream tasks. Previous methods typically focus on entropy-based exploration or empowerment-driven skill learning. However, entropy-based exploration struggles in large-scale state spaces (e.g., images), and empowerment-based methods with Mutual Information (MI) estimations have limitations in state exploration. To address these challenges, we propose a novel skill discovery objective that maximizes the deviation of the state density of one skill from the explored regions of other skills, encouraging inter-skill state diversity similar to the initial MI objective. For state-density estimation, we construct a novel conditional autoencoder with soft modularization for different skill policies in high-dimensional space. To incentivize intra-skill exploration, we formulate an intrinsic reward based on the learned autoencoder that resembles count-based exploration in a compact latent space. Through extensive experiments in challenging state and image-based tasks, we find our method learns meaningful skills and achieves superior performance in various downstream tasks.

025 026

027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 Reinforcement Learning (RL) has achieved remarkable success in game AI (Silver et al., 2018; Ye 029 et al., 2021), autonomous cars (Cao et al., 2023; Wu et al., 2022), and embodied agents (Hansen et al., 2022; Miki et al., 2022). Traditionally, RL agents rely on well-designed reward functions 031 to learn specific tasks (Luo et al., 2023). However, designing these reward functions is resource-032 intensive and often requires domain-specific expertise (Kwon et al., 2023; Gu et al., 2023), making 033 the learned policies dependent on handcrafted rewards and potentially unable to capture the com-034 plexity of real-world scenarios. This reliance limits the agent's generalization capability across diverse tasks and results in poor adaptability. In contrast, recent advances in Large Language Models (LLMs) (Han et al., 2021; Achiam et al., 2023) signify that unsupervised auto-regression has 036 led to powerful pre-trained language models, which can be adapted to downstream tasks via super-037 vised fine-tuning (Ouyang et al., 2022; Touvron et al., 2023). A powerful vision encoder can also be pre-trained via masked prediction without annotations or labels (He et al., 2022; Bardes et al., 2024; Grill et al., 2020), and the encoder can be used to solve various vision tasks (Majumdar et al., 040 2023; Nair et al., 2023). Inspired by these breakthroughs, it is desirable to further explore similar 041 unsupervised learning methods within the RL field. The goal is for unsupervised RL to learn useful 042 behaviors in the absence of external rewards, thus equipping them with the capacity to quickly adapt 043 to new tasks with limited interactions (Laskin et al., 2021).

044 The formulation of unsupervised RL has been studied in many prior works, which can be roughly categorized into empowerment-based skill discovery (Gregor et al., 2016) and pure exploration 046 methods (Liu & Abbeel, 2021b). Empowerment-based methods aim to maximize the Mutual In-047 formation (MI) between states and skills, and the MI term can be estimated by different variational 048 estimators (Song & Ermon, 2020). These methods have shown effectiveness in learning discriminative skills for state-based locomotion tasks (Eysenbach et al., 2019). However, the learned skills often have limited state coverage due to the inherent sub-optimality in the MI objective (Yang et al., 051 2023), which can lead to sub-optimal adaptation performance in downstream tasks and becomes more severe in large-scale state space (Park et al., 2024). Recent works introduce additional tech-052 niques like Lipschitz constraints and metric-aware abstraction to enhance the exploration abilities (Park et al., 2022; 2023; 2024). Pure exploration methods encourage the agent to explore the environment with maximum state coverage; however, this can lead to extremely dynamic skills rather
than meaningful behaviors for downstream tasks (Liu & Abbeel, 2021b; Laskin et al., 2022). Meanwhile, both the MI estimator and entropy estimation are not directly scalable to large-scale spaces,
such as pixel-based environments (Rajeswar et al., 2023; Park et al., 2024).

058 To overcome the aforementioned limitations, this work proposes a novel skill discovery method by maximizing the State Density Deviation of Different skills (SD3). Specifically, we construct a con-060 ditional autoencoder for state density estimation of different skills in high-dimensional state spaces. 061 Each skill policy is then encouraged to explore regions that deviate significantly from the state den-062 sity of other skills, which encourages inter-skill diversity and leads to discriminative skills. For a 063 stable state-density estimation of significantly different skills, we adopt soft modularization for the 064 conditional autoencoder to make the skill-conditional network a weighted combination of the shared modules according to a routing network determined by the skill. We show the skill-deviation ob-065 jective of SD3 resembles the initial MI objective in a special case. Further, to incentivize intra-skill 066 *exploration*, we formulate an intrinsic reward from the autoencoder based on the learned latent space, 067 which extracts the skill-relevant information and is scalable to large-scale problems. Theoretically, 068 such an intrinsic reward is closely related to the provably efficient count-based exploration in tabular 069 cases. To summarize, SD3 encourages inter-skill diversity via density deviation and intra-skill exploration via count-based exploration in a unified framework. We conduct extensive experiments in 071 Maze, state-based Unsupervised Reinforcement Benchmark (URLB), and challenging image-based 072 URLB environments, showing that SD3 learns exploratory and diverse skills. 073

Our contribution can be summarized as follows. (i) We propose a novel skill discovery objective based on state density deviation of skills, providing a straightforward way to learn diverse skills with different state occupancy. (ii) We propose a novel conditional autoencoder with soft modularization to estimate the state density of significantly different skills stably. (iii) The learned latent space of the autoencoder provides an intrinsic reward to encourage intra-skill exploration that resembles countbased exploration in tabular MDPs. (iv) Our method achieves state-of-the-art performance in various downstream tasks in challenging URLB benchmarks and demonstrates scalability in image-based URLB tasks. The open-sourced code is available at https://github.com/s7p77/SD3.

081 082

083

2 PRELIMINARIES

084 Markov Decision Process A Markov Decision Process (MDP) constitutes a foundational model 085 in decision-making scenarios. We consider the process of an agent interacting with the environment as an MDP with discrete skills, defined by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{Z}, \mathcal{P}, r, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is 087 the action space, \mathcal{Z} is the skill space, $\mathcal{P}: \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$ is the transition function, $r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ 088 is the reward function, and γ is the discount factor. In this work, we consider a discrete skill space 089 \mathcal{Z} that contains n skills since calculating the skill density deviation requires density estimation of 090 all skills, while SD3 can also be extended to a continuous skill space by sampling skills from a 091 continuous distribution for approximation. In each timestep, an agent follows a skill-conditional 092 policy $\pi(a|s, z)$ to interact with the environment. Given clear contexts, we refer to 'skill-conditional policy' as 'skill'.

094

Unsupervised RL Unsupervised RL typically contains two stages: unsupervised pre-training and 096 fast policy adaptation. In the unsupervised training stage, the agent interacts with the environment without any extrinsic reward. The policy $\pi(a|s, z)$ is learned to maximize some intrinsic rewards r_t formulated by an estimation of the MI term or the state entropy. The aim of unsupervised pre-098 training is to learn a set of useful skills that potentially solve various downstream tasks via fast policy adaptation. In the adaptation stage, the policy $\pi(a|s, z^{\star})$ with a chosen skill z^{\star} is optimized by RL 100 algorithms with certain extrinsic rewards to adapt to specific downstream tasks. In the following, 101 we denote $I(\cdot; \cdot)$ by the MI between two random variables and $\mathcal{H}(\cdot)$ by either the Shannon entropy 102 or differential entropy, depending on the context. We use uppercase letters for random variables and 103 lowercase letters for their realizations. We denote $d^{\pi}(s) \triangleq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi)$ as the 104 normalized probability that a policy π encounters state s. 105

The empowerment-based skill discovery algorithms try to estimate the MI between S and Z via $I(S; Z) = \mathbb{E}_{z \sim p(z), s \sim p^{\pi}(s|z)}[\log p(z|s) - \log p(z)]$. Given the computational challenges associated with the posterior p(z|s), a learned skill discriminator $q_{\phi}(z|s)$ is employed (Eysenbach et al., 2019) and a variational lower bound is established for the MI term as $I(Z;S) \geq \mathbb{E}_{z \sim p(z), s \sim p^{\pi}(s|z)} [\log q_{\phi}(z|s) - \log p(z)]$. Alternatively, pure exploration methods estimate state entropy by summing the log-distances between each particle and its k-th nearest neighbor, as $\mathcal{H}(s) \propto \sum_{s_i} \ln \|s_i - NN_k(s_i)\|$.

3 Method

117

118 119

120

In this section, we first introduce the proposed SD3 algorithm that performs skill discovery by maximizing inter-skill diversity via state density estimation. Next, we present the formulation of intrinsic rewards for intra-skill exploration. Finally, we provide a qualitative analysis of SD3.

3.1 SKILL DISCOVERY VIA DENSITY DEVIATION

121 We develop our skill discovery strategy from a straightforward intuition: The explored region of 122 each skill should deviate from other skills as far as possible. Formally, the optimizing objective for 123 skill discovery, denoted as I_{SD3} and referred to as *density deviation*, is defined by

 $I_{\text{SD3}} \triangleq \mathbb{E}_{z \sim p(z), s \sim d_z^{\pi}(s)} \left[\log \frac{\lambda \, d_z^{\pi}(s)}{\lambda \, d_z^{\pi}(s) p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s) p(z')} \right],$

124

125 126

127 where z is sampled from p(z), s is sampled from the state distribution induced by the skill policy 128 $\pi(a|s, z)$, and $\lambda > 0$ is a weight parameter. The numerator $d_z^{\pi}(\cdot)$ is the state density of skill z, and the 129 denominator is the weighted average of the state density of z and those of other skills $\{z'\}$. Since we 130 uniformly sample skills from the skill set that contains n skills, we have p(z) = 1/n for each skill z. 131 According to Eq. (1), it is easy to check that I_{SD3} attains its maximum when $\sum_{z'\neq z} d_{z'}^{\pi}(s) \to 0$ for 132 all (s, z) such that $p(z) \cdot d_z^{\pi}(s) > 0$, and the maximum value is $\mathcal{H}(Z)$. In this case, the state $s \sim d_z^{\pi}(\cdot)$ 133 visited by skill z has zero visitation probability by other skills, which means the explored regions of all skills do not overlap, and the learned skills are fully distinguishable. However, enforcing 134 such a strong objective to separate the overlapping explored areas of skills may lead to limited state 135 coverage for each skill. In extreme cases, each skill might only visit a distinct state that other skills 136 do not access. Although this leads to distinguishable skills, the overall state coverage becomes 137 overly limited, making them undesirable for learning meaningful behaviors. 138

In SD3, we adopt two mechanisms for addressing this problem. (i) A weight parameter λ is used in the learning objective to regularize the gradients of I_{SD3} to other skills. To see this, for each (s, z), we denote the state density of other skills $\{z'\}$ except z as $\rho_{z^c} \triangleq \sum_{z' \neq z} d_{z'}^{\pi}(s)$, then the gradient of $I_{SD3}(s, z)$ to ρ_{z^c} becomes

$$\nabla_{\rho_{z^c}} I_{\text{SD3}}(s, z) = -1/(\lambda d_z^{\pi}(s) + \rho_{z^c}(s)), \tag{2}$$

(1)

where $I_{\text{SD3}}(s, z)$ is the density ratio for a specific (s, z) and the proof is attached in A.1. Thus, for skill z, increasing λ will weaken the gradient of SD3 in reducing the state densities of other skills, which prevents skill collapse in SD3. (ii) We introduce explicit intra-skill exploration based on the latent space learned in estimating the skill density, which will be discussed in §3.2. To maximize I_{SD3} , we adopt a modified Conditional Variational Auto-Encoder (CVAE) to stably estimate the state density for skills, which we introduce as follows.

$$\log d_{z}^{\pi}(s) = \mathbb{E}_{Q(h|s,z)} \log \left[P(s|z) \right] = \mathbb{E}_{Q(h|s,z)} \log \left[\frac{P(s,h|z)}{Q(h|s,z)} \right] + \mathbb{E}_{Q(h|s,z)} \log \left[\frac{Q(h|s,z)}{P(h|s,z)} \right]$$
$$\geq \mathbb{E}_{Q(h|s,z)} \log \left[\frac{P(s|h,z)P(h|z)}{Q(h|s,z)} \right] = \underbrace{\mathbb{E}_{Q(h|s,z)} \log \left[P(s|h,z) \right] - D_{\mathrm{KL}} \left[Q(h|s,z) \| P(h|z) \right]}_{\mathcal{L}_{z}^{\mathrm{elbo}}(s)}, \tag{3}$$

159 160

156 157

143 144

where the latent vector h is sampled from a variational posterior distribution (i.e., Q(h|s, z)) conditioned on the state and skill, and the inequality holds by dropping off the non-negative second



Figure 1: An overview of the CVAE architecture. (a) The encoder-decoder network with soft modularization. The feature extractor of state can be MLPs or convolution layers according to state- or image-based environment. (b) The inter-skill diversity objective for skill discovery and the intraskill intrinsic reward for exploration can be derived from the learned CVAE.

179 term, which is the definition of $D_{\rm KL}(Q(h|s,z)||P(h|s,z))$. Meanwhile, we use P(s,h|z) =180 P(h|z)P(s|h,z) to decompose the joint distribution. According to Eq. (3), maximizing the Evi-181 dence Luwer-Bound (ELBO) $\mathcal{L}_z^{\text{elbo}}(s)$ can approximate the skill-conditioned state distribution, as $\log d_z^{\pi}(s) \approx \max_Q \mathcal{L}_z^{\text{elbo}}(s)$. To maximize $\mathcal{L}_z^{\text{elbo}}(s)$, we learn an encoder network $Q_{\phi}(h|s,z)$ to ob-182 183 tain the posterior of latent representation, where the posterior is represented by a diagonal Gaussian. Then, a latent vector h is sampled from the posterior, and a decoder network $P_{\psi}(s|h,z)$ is used to 185 reconstruct the state. The KL-divergence in $\mathcal{L}_z^{\text{elbo}}(s)$ regularizes the latent space via a prior distri-186 bution P(h|z), which is set to a standard Gaussian. The whole objective is optimized via stochastic 187 gradient ascent with a reparameterization trick (Kingma & Welling, 2013; Kingma et al., 2019). To calculate I_{SD3} , we perform state density estimations for all skills via forward inference based on 188 the learned encoder and decoder. In calculating I_{SD3} , we adopt efficient parallelization to calculate 189 $\mathcal{L}_z^{\text{elbo}}(s)$ for all skills $z \in \mathcal{Z}$ in one forward pass, which minimizes the run-time increase with the 190 number of skills. 191

192

205 206

175

176

177

178

193 **Soft Modularization for CVAE** As we maximize the state-density deviation in skill discovery, the resulting skills become diverse, and the corresponding state occupancy for different skills tends 194 to be very different. In CVAE-based density estimation, since different skills share the same network 195 parameters, optimizing $\mathcal{L}_{z}^{\text{elbo}}$ for one skill can negatively affect the density estimation of other skills 196 with significantly different state densities. Empirically, we also find obtaining an accurate estimation 197 of $d_z^{\pi}(s)$ for all skills $z \in \mathcal{Z}$ can be difficult. As a result, we adopt a soft modularization technique that automatically generates soft network module combinations for different skills without explicitly 199 specifying structures. As shown in Figure 1, the soft modularized CVAE contains an unconditional 200 basic network and a routing network, where the routing network takes the skill and state embedding 201 as input to estimate the routing strategy. Suppose each layer of the encoder/decoder network has 202 m modules, then the routing network gives the probabilities $p \in \mathbb{R}^{m \times m}$ to weight modules con-203 tributing to the next layer. Specifically, considering *l*-th layer has probabilities $p^l \in \mathbb{R}^{m \times m}$, then the 204 probability in the next layer is

$$p^{l+1} = \mathcal{W}^l \big(\operatorname{ReLU}(g(p^l) \odot (u \odot v)) \big), u = f_1(s), v = f_2(z), \tag{4}$$

207 where \odot denotes element-wise product, $g(\cdot)$, $f_1(\cdot)$ and $f_2(\cdot)$ are all fully connected layers that $f_1(\cdot)$ 208 and $f_2(\cdot)$ map state s and skill z to the same dimensions (e.g., d), and $g(\cdot)$ maps p^l to the dimension 209 d. Then we have $\mathcal{W}^l \in \mathbb{R}^{m^2 \times d}$ to project the joint feature to a probability vector of layer l+1. In the 210 basic network, we denote the input feature for the j-th module in the l-the layer as $g_i^l \in \mathbb{R}^d$; then we 211 have $g_i^{l+1} = \sum_j \hat{p}_{i,j}^l(\text{ReLU}(\mathcal{W}_j^l g_j^l))$ for the next layer, where $\hat{p}_{i,j}^l = \exp(p_{i,j}^l)/(\sum_{j=1}^m \exp(p_{i,j}^l))$ 212 is the normalized vector that weights the j-th module in the l-th layer to contribute to the i-th module 213 in the l + 1-th layer. We remark that the soft modularization technique was originally proposed in 214 multi-task RL (Yang et al., 2020), while we extend it to encoder-decoder-based CVAE for density 215 estimation. The detailed architecture is given in §B.2.



Figure 2: An illustration of skill discovery in SD3. The skills start with overlapping areas and are separated via state-density deviation. Then, each skill explores the environment independently, resulting in overlapped but expanded areas. SD3 separates the areas again and leads to distinguishable skills. Such a process repeats and ultimately leads to exploratory and diverse skills.

233 3.2 LATENT SPACE EXPLORATION

228

229

230

231

232

234

241 242

250

264

As we discussed above, the SD3 objective that only maximizes the density deviation may lead to skill collapse. In addition to introducing an additional parameter λ in Eq. (1), we find the learned CVAE in Figure 1 can provide a *free-lunch* intrinsic reward for efficient intra-skill exploration. In SD3, we derive an intrinsic reward based on the latent space that learns skill-conditioned representations for states. Specifically, the KL-divergence term $D_{\text{KL}}[Q(h|s, z)||r(h)]$ in CVAE objective serves as an upper bound of the conditional MI term I(S; H|Z), as

$$I(S; H|Z) = \mathbb{E}_{p(s,z), Q_{\phi}(h|s,z)} \Big[\log Q_{\phi}(h|s,z) / P(h|z) \Big] \le \mathbb{E}_{p(s,z), Q_{\phi}(h|s,z)} \Big[\log Q_{\phi}(h|s,z) / r(h) \Big],$$
(5)

where *H* denotes the random variable of the sampled latent representation *h*, and r(h) the prior distribution set to a standard Gaussian, and $P(h|z) \triangleq \mathbb{E}_{P(s|z)}Q_{\phi}(h|s,z)$. The inequality holds since $D_{\mathrm{KL}}[P(h|z)||r(h)] \ge 0$ for all $z \in \mathbb{Z}$. Since $D_{\mathrm{KL}}[Q_{\phi}(h|s,z)||r(h)]$ is constrained in CVAE learning, the MI between states and latent representations for each skill is also compressed according to Eq. (5). Thus, the latent space in CVAE learns a compressive representation while retaining important information as the representation is then used for reconstruction. Based on the learned representation, we define the intrinsic reward for intra-skill exploration as

$$r_z^{\exp}(s) = D_{\mathrm{KL}}[Q_{\phi}(h|s,z) \| r(h)],$$
 (6)

251 where $Q_{\phi}(h|s,z)$ is the posterior network learned in CVAE. The intrinsic reward in Eq. (6) quantifies the degree of compression of representation with respect to the state, which measures skill-253 conditioned state novelty in a compact space for intra-skill exploration. Intuitively, if a state $s^{(1)}$ is 254 frequently visited by skill z, then the corresponding latent distribution is close to r(h) according to 255 Eq. (5), and the resulting reward $r_z^{exp}(s^{(1)})$ will be close to zero. In contrast, if a state $s^{(2)}$ is novel for 256 skill z, then the corresponding intrinsic reward will be high since the latent posterior $Q_{\phi}(h|s^{(2)},z)$ 257 can be very different from the prior r(h). Thus, in exploration, such reward encourages the policy to 258 find the scarcely visited states $\{s^+\}$ (with a high $D_{\mathrm{KL}}[Q_{\phi}(h|s,z)||r(h)]$) and explore these states. 259

An illustration of the skill learning process of SD3 is shown in Figure 2. The state occupancy of different skills overlaps initially in Figure 2(i), then we maximize I_{SD3} via per-instance estimation and set it to an intrinsic reward as

$$r_{z}^{\text{sd3}}(s) = \log \frac{\lambda \, d_{z}^{\pi}(s)}{\lambda \, d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')},\tag{7}$$

which encourages skill density deviation and leads to more diverse skills with separate state coverage, as in Figure 2(ii). Then the exploration reward $r_z^{\exp}(s)$ is used to encourage intra-skill exploration, which makes each skill explore unknown areas independently. After exploration, the state coverage of each skill increases and may lead to state-coverage overlapping again among skills, as in Figure 2(iii). Then the density-derivation reward $r_z^{\operatorname{sd3}}(s)$ re-separates the updated areas to obtain distinguished skills, as in Figure 2(iv). The above process repeats for many rounds and SD3 finally learns exploratory and diverse skills. The algorithmic description of our method is given in
 Algorithm 1.

2732743.3 QUALITATIVE ANALYSIS

In this section, we give a qualitative analysis of the proposed SD3 objective and exploration reward, which encourage inter-skill diversity and intra-skill exploration, respectively.

The skill discovery objective I_{SD3} in Eq. (1) leads to diverse skills with separate explored areas, which is similar to the MI-based skill discovery objectives. As we usually set $\lambda \ge 1$ to prevent skill collapse, the following theorem connects I_{SD3} and the previous MI objectives.

Theorem 3.1. With $\lambda \ge 1$, we have

$$I(S;Z) \le I_{SD3} \le c_0 + I(S;Z).$$
 (8)

where $c_0 = \log \lambda$. Specially, $I_{SD3} = I(S; Z)$ if $\lambda = 1$.

The above theorem shows when we maximize skill deviation via I_{SD3} , the MI between S and Z also increases. The previous MI objective becomes a special case of I_{SD3} , where the introduced λ provides flexibility to control the strength of skill deviation. In the following, we connect the proposed intrinsic reward to the provably efficient count-based exploration in tabular cases.

Note that since λ only relates to the overall objective I_{SD3} and does not affect the estimation of state density, the exploration bonus holds for arbitrary $\lambda \ge 1$.

Theorem 3.2. In tabular MDPs, optimizing the intra-skill exploration reward is equivalent to countbased exploration, as

298

299

300

301

292

282

283

284

 $r_z^{\exp}(s) \approx \frac{|\mathcal{S}|/2}{N(s,z) + \kappa}.$ (9)

where N(s, z) is the count of visitation of state-skill pair (s, z) in experiences, |S| is the total number of states in a tabular case, and $\kappa > 0$ is a small non-negative constant.

As a result, maximizing the intra-skill exploration reward is equivalent to performing count-based exploration in previous works (Kolter & Ng, 2009; Strehl & Littman, 2008), which is provable efficient in tabular MDPs (Bellemare et al., 2016; Ostrovski et al., 2017). Through the approximation in a compact latent space, the intra-skill exploration encourages skill-conditional policy to increase the pseudo-count of rarely visited state-skill pairs in a high-dimensional space.

4 RELATED WORK

306 **Unsupervised Skill Discovery** Unsupervised skill discovery in RL aims to acquire a repertoire of 307 useful skills without relying on extrinsic rewards. Early efforts, such as VIC (Gregor et al., 2016), 308 DIAYN (Eysenbach et al., 2019), and DADS (Sharma et al., 2020), maximize the MI between the 309 skill and the state to discover diverse skills. However, as noted in EDL (Campos et al., 2020), LSD (Park et al., 2022), and CSD (Park et al., 2023), such MI-based methods usually prefer static 310 skills caused by poor state coverage and may hinder the application for downstream tasks. Recent 311 methods strive to address this limitation to learn dynamic and meaningful skills. These methods 312 perform explicit exploration or enforce Lipschitz constraints in the representation to maximize the 313 traveled distances of skills. Further, CIC (Laskin et al., 2022) employs contrastive learning between 314 state transitions and skills to encourage agent's diverse behaviors. BeCL (Yang et al., 2023) uses 315 contrastive learning to differentiate between various behavioral patterns and maximize the entropy 316 implicitly. ReST (Jiang et al., 2022) encourages the trained skill to stay away from the estimated state 317 visitation distributions of other skills. Some methods, like DISCO-DANCE (Kim et al., 2023), APS 318 (Liu & Abbeel, 2021a), SMM (Lee et al., 2020) and DISDAIN (Strouse et al., 2022), focus on in-319 troducing an auxiliary exploration reward to address insufficient exploration. Furthermore, to verify 320 the effectiveness of skill discovery in large-scale state space (e.g., images), recent methods including 321 Choreographer (Mazzaglia et al., 2023) and Metra (Park et al., 2024) evaluate the effectiveness of methods on pixel-based URLB (Rajeswar et al., 2023), which often relies on model-based agents to 322 learn meaningful knowledge from imagination, and skills are discovered in the latent space. Metra 323 (Park et al., 2024) constructs a latent space associated with the original state space via a temporal



Figure 3: Results for maze experiment. We visually demonstrate the agent's ability to explore the environment and the diversity of skills discovered by the agent. The agent starts from the black dot of the maze and interacts for 250K steps. Both DIAYN and DADS do not reach the right side of the maze while obtaining distinguishable trajectories highlighted by different colors. The trajectories of CIC span the entire maze but appear chaotic. In contrast, SD3 can reach the farthest position from the starting point and facilitates easy differentiation of trajectories of different skills.

distance metric, which enables skill learning in high-dimensional environments by maximizing the coverage. In contrast, our method promotes skill diversity by encouraging deviations in skill density and enhances state coverage through latent space exploration. We validate our approach's efficacy through experiments on state-based and pixel-based tasks across various environments.

Unsupervised RL According to URLB (Laskin et al., 2021), URL algorithms are classified into three main categories: knowledge-based, data-based, and competence-based. Knowledge-based al-gorithms (Pathak et al., 2017; 2019; Burda et al., 2019) leverage the agent's predictive capacity or understanding of the environment, and the intrinsic reward is tied to the novelty of the agent's behaviors, encouraging the agent to explore areas where its model is less certain. Data-based al-gorithms (Liu & Abbeel, 2021b; Yarats et al., 2021) maximize the state entropy to maximize state coverage of skills. Competence-based algorithms (Lee et al., 2020; Eysenbach et al., 2019; Liu & Abbeel, 2021a; Nieto et al., 2021) pre-train the agent to learn useful skills that can be utilized to complete downstream tasks. Our method can be categorized as competence-based, while also combining the benefit of knowledge-based algorithms to encourage exploration. In addition, some recent algorithms do not easily fit into these categories. For example, LCSD (Ju et al., 2024) estab-lishes connections between skills, states, and linguistic instructions to guide task completion based on external language directives. DuSkill (Kim et al., 2024) utilizes a guided diffusion model to gen-erate versatile skills beyond dataset limitations, thereby enhancing the robustness of policy learning across diverse domains. EUCLID (Yuan et al., 2023) improves downstream policy learning per-formance by jointly pre-training dynamic models and unsupervised exploration strategies. VGCRL (Choi et al., 2021) applies variational empowerment to learn effective state representations, thereby improving exploration.

5 EXPERIMENTS

We start by introducing experiments in Maze to visualize the skills. Subsequently, we validate the effectiveness of SD3 by conducting experiments on challenging tasks from the DeepMind Control Suite (DMC) (Tassa et al., 2018), with both state-based (Laskin et al., 2021) and pixel-based (Rajeswar et al., 2023) observations. Finally, we conduct ablation studies to demonstrate the factors that influence the effectiveness of SD3.

371 5.1 MAZE EXPERIMENT

We conduct experiments in a 2D maze to visually demonstrate the learned skills, as shown in Figure 3. The agent's initial state is represented by a black dot, with different colored lines indicating the trajectories corresponding to the different skills it has learned. The agent's state is the current positional information, and the actions represent the velocity and direction of movement. Building on this, we compare SD3 with two classical MI-based methods, DIAYN (Eysenbach et al., 2019) and DADS (Sharma et al., 2020), whose objectives correspond to the reverse form $\mathcal{H}(Z) - \mathcal{H}(Z|S)$ and the forward form $\mathcal{H}(S) - \mathcal{H}(S|Z)$ of the MI term I(S; Z), respectively. Additionally, we com-



Figure 4: Results for state-based URLB. The aggregate statistics (Agarwal et al., 2021) indicate the adaptation performance of different unsupervised RL methods in 12 downstream tasks. In terms of IQM, Mean, and OG metrics, SD3 outperforms other competence-based methods and significantly surpasses pure exploration methods, achieving 77.37%, 76.19%, and 23.91%, respectively.

pare SD3 with an entropy-based CIC algorithm (Laskin et al., 2022), whose primary objective is to maximize state-transition entropy $\mathcal{H}(\tau)$ to generate diverse behaviors. We employ the PPO as the backbone and train n = 10 skills for each algorithm.

396 We delineate the learned skills of each algorithm within the maze environment in Figure 3 and intro-397 duce two key metrics for comparing SD3 with other methods: state coverage and distinguishability 398 of skills, where insufficient state coverage may impede the acquisition of dynamic skills, and the 399 lack of distinguishability leads to similar behaviors of skills. According to the results, (i) DIAYN 400 and DADS fail to extend to the upper-right corner of the maze, but exhibit clear distinctions among 401 trajectories of skills, indicating that merely maximizing I(S; Z) can learn discriminable skills but 402 lack effective exploration of the state space; (ii) CIC demonstrates the best state coverage while learns skills with mixed trajectories due to the maximization of $\mathcal{H}(s)$ as its primary objective; (iii) 403 In contrast, SD3 strikes a balance between state coverage and empowerment in skill discovery. It 404 learns discriminable skills by maximizing the deviation between the state densities of a certain skill 405 and others. Meanwhile, SD3 achieves commendable state coverage through latent space exploration. 406

408 5.2 STATE-BASED URLB

388

389

390

391

392

407

409 According to state-based URLB (Laskin et al., 2021), we evaluate our approaches in 12 downstream 410 tasks across 3 distinct continuous control domains, each designed to evaluate the effectiveness of 411 algorithms under high-dimensional state spaces. The three domains are Walker, Quadruped, and 412 Jaco Arm. Specifically, Walker involves a biped constrained to a 2D vertical plane with a state 413 space $S \in \mathbb{R}^{24}$ and an action space $\mathcal{A} \in \mathbb{R}^{6}$. The agent in the *Walker* domain must learn to 414 maintain balance and move forward, completing four downstream tasks: stand, walk, run, and flip. 415 *Quadruped* features a four-legged robot in a 3D environment, characterized by a state space $S \in \mathbb{R}^{78}$ and an action space $\mathcal{A} \in \mathbb{R}^{16}$. The downstream tasks, including *stand*, *run*, *jump*, and *walk*, pose 416 417 challenges to the agent due to the complex dynamics of its movements. Jaco employs a 6-DOF robotic arm with a three-finger gripper, functioning within a state space $S \in \mathbb{R}^{55}$ and an action 418 state $\mathcal{A} \in \mathbb{R}^9$. Primary downstream tasks in *Jaco* Arm include reaching and manipulating objects at 419 various positions. 420

421 Baselines. We conduct comparisons between SD3 and the baselines delineated across the three 422 URL algorithm categories as defined by URLB (Laskin et al., 2021). These categories encompass 423 knowledge-based baselines, which consist of ICM (Pathak et al., 2017), Disagreement (Pathak et al., 2019), and RND (Burda et al., 2019); data-based baselines, which include APT (Liu & Abbeel, 424 2021b) and ProtoRL (Yarats et al., 2021); and competence-based baselines, comprising SMM (Lee 425 et al., 2020), DIAYN (Eysenbach et al., 2019), and APS (Liu & Abbeel, 2021a). Furthermore, we 426 extend our comparisons to include other novel competence-based algorithms such as CSD (Park 427 et al., 2023), Metra (Park et al., 2024), BeCL (Yang et al., 2023), and CIC (Laskin et al., 2022). 428

Evaluation. We employ a rigorous evaluation to assess the performance of SD3 alongside other algorithms, involving a two-phase process. Initially, a pre-training of 2M steps is performed using only intrinsic rewards, followed by a fine-tuning phase of 100K steps on each downstream task using extrinsic rewards. Building upon prior work (Laskin et al., 2021), we utilize DDPG as the backbone



Figure 5: (a) We conduct experiments on pixel-based URLB to demonstrate the scalability of SD3 for large-scale problems. (b) It can be observed that SD3 retains higher performance ratio than CIC in the noisy domain.

449 algorithm. To ensure statistical rigor and mitigate the impact of incidental factors in RL training, we conduct experiments across multiple seeds (10 seeds per algorithm), resulting in a substantial 450 volume of experimental runs (i.e., 1560 = 13 algorithms $\times 10$ seeds $\times 3$ domains $\times 4$ tasks). We 451 employ four statistical metrics to assess performance: Median, interquatile mean (IQM), Mean, and 452 optimality gap (OG) (Agarwal et al., 2021). IQM focuses on the central tendency of the middle 453 50%, excluding the top and bottom quartiles. OG understands the extent to which the algorithm 454 approaches the optimal level, where the optimal level is determined by the expert models' ultimate 455 score obtained on each downstream task. 456

Results. According to Figure 4, SD3 achieves the highest IQM score at 77.37%, slightly surpassing 457 CIC and BeCL, which scores 75.19% and 75.38% respectively, and significantly outperforming 458 other competence-based algorithms such as Metra (61.01%), CSD (54.93%), and APS (43.61%). On 459 the OG metric, SD3's gap to optimal performance is 23.91%, marginally better than CIC and BeCL 460 at 25.65% and 25.44%, respectively, and far superior to Metra (39.25%), CSD (42.43%), and APS 461 (55.76%). Additionally, compared to purely exploratory methods, SD3 significantly outperforms 462 the best-performing method, APT, on both IQM and OG metrics, with APT scoring 67.74% and 463 34.98% on these metrics, respectively. The remarkable performance of SD3 stems from two main factors. First, the use of r^{sd3} facilitates the learning of distinguishable skills by the agent, thereby 464 465 facilitating effective adaptation across various downstream tasks. Second, the learned compressed 466 representation of the high-dimensional state space leads to efficient intra-skill exploration within a compact space, which not only maintains skill consistency but also enhances exploration ability. 467

468

445

446

447 448

469 5.3 PIXEL-BASED URLB 470

To further validate the effectiveness of SD3, we conduct experiments on pixel-based URLB (Rajeswar et al., 2023), which includes *Walker* and *Quadruped* domains with 8 downstream tasks. The pixel-based environment employs raw pixel data as input, foregoing abstracted features, or processed sensor information. The challenge of deriving meaningful skills from such unrefined inputs is substantial, particularly in the absence of external rewards. Meanwhile, exploration becomes more difficult in image-based spaces, thereby testing the exploration ability of algorithms under conditions that closely resemble practical applications.

Baselines. We compared SD3 with the top three performing algorithms in state-based experiments, i.e., BeCL (Yang et al., 2023), CIC (Laskin et al., 2022), and APT (Liu & Abbeel, 2021b), as well as with the recently proposed skill discovery algorithms including CSD (Park et al., 2023) and Metra (Park et al., 2024). Among these, APT stands out as a data-based algorithm, which can also be considered a representative of pure exploration algorithms and demonstrates strong performance in exploring environments. The others are competence-based algorithms, which accomplish downstream tasks by learning useful and diverse skills.

Evaluation. We conduct 2M steps of pre-training solely based on intrinsic rewards in each domain, followed by 100K steps of fine-tuning on the downstream tasks using extrinsic rewards. The

486 scores achieved in the downstream tasks are used to evaluate the algorithm. According to the offi-487 cial benchmark of the pixel-based URLB (Rajeswar et al., 2023), unsupervised RL algorithms often 488 perform poorly when combined with a model-free method (e.g., DDPG (Lillicrap et al., 2016) or 489 DrQv2 (Yarats et al., 2022)) with image observations, while performing much better when using a 490 model-based backbone (e.g., Dreamer (Hafner et al., 2021)). Thus, we follow this setting and conduct experiments with Dreamer backbone. We report the average adaptation performance in Figure 491 5(a). In the relatively simple *Walker* domain, SD3 achieves the best performance (93.42%), slightly 492 outperforming other methods (i.e., CIC-91.29%, APT-88.17%, CSD-84.26%). In the challenging 493 Quadruped domain, SD3 outperforms CIC (77.57% and 75.89%, respectively) and shows signif-494 icant improvement over other competence-based methods (i.e., CSD-65.89%, Metra-53.53%) and 495 the best pure-exploration method in state-based URLB (i.e., APT-61.96%). This highlights SD3's 496 commendable advantages in both various image-based tasks. 497

498 499

500

5.4 ROBUSTNESS EXPERIMENT

501 Unlike CIC, APS, and BeCL, which rely on entropy-based exploration strategies, SD3 introduces 502 a novel exploration reward that resembles a UCB-style bonus. Such a UCB-term in exploration 503 is provable efficient in linear and tabular MDPs, which has been rigorously studied in previous 504 research (Jin et al., 2023; ZHANG et al., 2021). In contrast, the entropy-based exploration used 505 in previous methods has the disadvantage of being non-robust (e.g., adding small noise will signifi-506 cantly affect its entropy). Thus, to further verify that the robustness of SD3, we conduct experiments 507 in noisy domains of URLB by adding noise during pre-training, which is sampled from N(0, 0.1), 508 followed by noise-free fine-tuning to assess the learned skills.

Evaluation. We choose CIC for comparison, which performs competitively with our method in standard URLB. Each technique is evaluated across 5 random seeds and the results are given in Figure 5(b). The Performance Ratio (PR) denotes the ratio of the adaptation score in the noisy domain to that in the normal setup. According to the results, it is evident that the UCB-bonus used in SD3 is more robust than entropy-based rewards in noisy environments, achieving significantly higher Performance Ratio than CIC. The detailed results are attached in Appendix E.

514 515

517

516 5.5 Ablation Studies and Visualization

518 We provide ablation studies for components in skill discovery and skill adaptation of SD3. For skill discovery, we perform the comparison on (i) density estimation with and without soft modulariza-519 tion, and (ii) the different settings of temperatures in the routing network. The final rewards for skill 520 discovery contain $r_z^{sd3}(s)$ and $r_z^{exp}(s)$. We conduct ablation studies on (iii) different settings of λ in 521 calculating $r_z^{sd3}(s)$, as well as (iv) the different balance factors of the two rewards. For skill adap-522 tation, we sampled skills randomly to evaluate their generalization ability in our main results. In 523 ablation studies, (v) we evaluate two more skill-choosing strategies in adaptation for a comparison. 524 We refer to Appendix D for detailed results and analysis. We also provide visualizations of skills learned in tree-like Maze and DMC tasks in Appendix C. The results show that SD3 learns dynamic 526 and valuable skills, enabling the agent to adapt to downstream tasks quickly.

527 528

6 CONCLUSION

529 530

We propose a novel skill discovery method that promotes skill diversity by encouraging skill deviations in state density and enhancing state coverage through latent space exploration. We realize a novel soft modularization architecture for state density estimation of different skills. Theoretically, the skill discovery objective also maximizes the initial MI term, and the resulting intra-skill exploration bonus resembles count-based exploration. Moreover, our four experiments complement each other and collectively provide sufficient evidence that SD3 demonstrates superior and more comprehensive performance compared to other methods. One limitation of our method is that the soft modularization architecture is limited to discrete skill spaces, and the theoretical analysis of the exploration bonus requires the assumption of tabular MDPs. In the future, we will extend the idea of skill discovery to LLM-based agents to learn meaningful skills in more complex environments.

540 7 ETHICS STATEMENT

This work does not involve any human subjects or personally identifiable information. This work focuses on the field of unsupervised RL, and therefore does not involve any datasets. No sensitive data or unethical methodologies are employed. We declare no conflicts of interest related to the sponsorship or publication of this work. The research has adhered to the ICLR Code of Ethics, with special attention to fairness and bias concerns.

547 548 549

555

556

561

562

563

566

567

568 569

570

571

572

579

580

581 582

583

584

585

586

588

589

590

8 REPRODUCIBILITY

All experiments and results reported in this paper can be reproduced using the provided anony anonymous source code. Details regarding the model architecture and training parameters are included in Appendix B. The theorems discussed in section 3.3 are supported by detailed proofs provided in Appendix A.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
 - Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.
- Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational information
 bottleneck. In *International Conference on Learning Representations*, 2017.
 - Peter Auer and Ronald Ortner. Logarithmic online regret bounds for undiscounted reinforcement learning. Advances in neural information processing systems, 19, 2006.
 - Chenjia Bai, Lingxiao Wang, Lei Han, Animesh Garg, Jianye Hao, Peng Liu, and Zhaoran Wang. Dynamic bottleneck for robust self-supervised exploration. *Advances in Neural Information Processing Systems*, 34:17007–17020, 2021.
- Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mido
 Assran, and Nicolas Ballas. V-JEPA: Latent video prediction for visual representation learning,
 2024. URL https://openreview.net/forum?id=WFYbBOEOtv.
- Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
 Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
 - Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019.
 - Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. In *International Conference on Machine Learning*, pp. 1283–1294. PMLR, 2020.
 - Víctor Campos, Alexander Trott, Caiming Xiong, Richard Socher, Xavier Giró-i Nieto, and Jordi Torres. Explore, discover and learn: Unsupervised discovery of state-covering skills. In *International Conference on Machine Learning*, pp. 1317–1327. PMLR, 2020.
 - Zhong Cao, Kun Jiang, Weitao Zhou, Shaobing Xu, Huei Peng, and Diange Yang. Continuous improvement of self-driving cars using dynamic confidence-aware reinforcement learning. *Nature Machine Intelligence*, 5(2):145–158, 2023.
- Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Varia tional empowerment as representation learning for goal-based reinforcement learning. *CoRR*, abs/2106.01404, 2021. URL https://arxiv.org/abs/2106.01404.

- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.
- 598 Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv* 599 *preprint arXiv:1611.07507*, 2016.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena
 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
 et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone Tao, Xinyue Wei, Yunchao Yao, Xiaodi Yuan, Pengwei Xie, Zhiao Huang, Rui Chen, and Hao Su. Maniskill2: A unified benchmark for generalizable manipulation skills. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=b_CQDy9vrD1.
- Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In *ICLR*, 2021.
- Ku Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao,
 Ao Zhang, Liang Zhang, et al. Pre-trained models: Past, present and future. *AI Open*, 2:225–250,
 2021.
- ⁶¹⁵
 ⁶¹⁶ Nicklas A Hansen, Hao Su, and Xiaolong Wang. Temporal difference learning for model predictive control. In *International Conference on Machine Learning*, pp. 8387–8406. PMLR, 2022.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2015.
- Zheyuan Jiang, Jingyue Gao, and Jianyu Chen. Unsupervised skill discovery via recurrent skill train ing. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *Mathematics of Operations Research*, 48(3):1496–1521, 2023.
- Zhaoxun Ju, Chao Yang, Fuchun Sun, Hongbo Wang, and Yu Qiao. Rethinking mutual information for language conditioned skill discovery on imitation learning. In 34th International Conference on Automated Planning and Scheduling, 2024. URL https://openreview.net/forum? id=8VdptRkRYW.
- Hyunseung Kim, Byungkun Lee, Sejik Park, Hojoon Lee, Dongyoon Hwang, Kyushik Min, and
 Jaegul Choo. Learning to discover skills with guidance. In *Advances in Neural Information Processing Systems*, 2023.
- Woo Kyung Kim, Minjong Yoo, and Honguk Woo. Robust policy learning via offline skill diffusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13177–13184, 2024.
- 642 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends*® *in Machine Learning*, 12(4):307–392, 2019.
- ⁶⁴⁷ J Zico Kolter and Andrew Y Ng. Near-bayesian exploration in polynomial time. In *Proceedings of the 26th annual international conference on machine learning*, pp. 513–520, 2009.

648 649 650	Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with language models. In <i>The Eleventh International Conference on Learning Representations</i> , 2023.
651 652 653	Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang, Lerrel Pinto, and Pieter Abbeel. URLB: Unsupervised reinforcement learning benchmark. In <i>Neural Information Processing Systems (Datasets and Benchmarks Track)</i> , 2021.
654 655 656 657	Michael Laskin, Hao Liu, Xue Bin Peng, Denis Yarats, Aravind Rajeswaran, and Pieter Abbeel. Unsupervised reinforcement learning with contrastive intrinsic control. In <i>Advances in Neural</i> <i>Information Processing Systems</i> , 2022.
658 659 660	Lisa Lee, Benjain Eysenbach, Emilio Parisotto, Erix Xing, Sergey Levine, and Ruslan Salakhutdi- nov. Efficient exploration via state marginal matching, 2020. URL https://openreview. net/forum?id=Hkla1eHFvS.
661 662 663 664	Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In <i>ICLR</i> (<i>Poster</i>), 2016. URL http://arxiv.org/abs/1509.02971.
665 666	Hao Liu and Pieter Abbeel. Aps: Active pretraining with successor features. In International Conference on Machine Learning, pp. 6736–6747. PMLR, 2021a.
668 669	Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. In Advances in Neural Information Processing Systems, volume 34, pp. 18459–18473, 2021b.
670 671 672 673	Yongle Luo, Yuxin Wang, Kun Dong, Qiang Zhang, Erkang Cheng, Zhiyong Sun, and Bo Song. Relay hindsight experience replay: Self-guided continual reinforcement learning for sequential object manipulation tasks with sparse rewards. <i>Neurocomputing</i> , 557:126620, 2023.
674 675 676 677	Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. Where are we in the search for an artificial visual cortex for embodied intelligence? In <i>Advances in Neural Information Processing Systems</i> , volume 36, 2023.
678 679 680 681	Pietro Mazzaglia, Tim Verbelen, Bart Dhoedt, Alexandre Lacoste, and Sai Rajeswar. Choreographer: Learning and adapting skills in imagination. In <i>International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=PhkWyijGi5b.
682 683 684	Takahiro Miki, Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hut- ter. Learning robust perceptive locomotion for quadrupedal robots in the wild. <i>Science Robotics</i> , 7(62):eabk2822, 2022.
685 686 687 688	Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A univer- sal visual representation for robot manipulation. In <i>Conference on Robot Learning</i> , pp. 892–909. PMLR, 2023.
689 690 691	Juan José Nieto, Roger Creus, and Xavier Giro-i Nieto. Unsupervised skill-discovery and skill-learning in minecraft. <i>arXiv preprint arXiv:2107.08398</i> , 2021.
692 693 694	Georg Ostrovski, Marc G Bellemare, Aäron Oord, and Rémi Munos. Count-based exploration with neural density models. In <i>International conference on machine learning</i> , pp. 2721–2730. PMLR, 2017.
695 696 697 698 699	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35: 27730–27744, 2022.
700 701	Seohong Park, Jongwook Choi, Jaekyeom Kim, Honglak Lee, and Gunhee Kim. Lipschitz- constrained unsupervised skill discovery. In <i>International Conference on Learning Represen-</i> <i>tations</i> , 2022.

702 703 704	Seohong Park, Kimin Lee, Youngwoon Lee, and Pieter Abbeel. Controllability-aware unsupervised skill discovery. In <i>International Conference on Machine Learning</i> , volume 202, pp. 27225–27245, 2023.
705 706 707 708	Seohong Park, Oleh Rybkin, and Sergey Levine. METRA: Scalable unsupervised RL with metric- aware abstraction. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=c5pwL0Soay.
709 710 711	Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In <i>International Conference on Machine Learning</i> , pp. 2778–2787. PMLR, 2017.
712 713 714	Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement. In International Conference on Machine Learning, pp. 5062–5071. PMLR, 2019.
715 716 717	Sai Rajeswar, Pietro Mazzaglia, Tim Verbelen, Alexandre Piché, Bart Dhoedt, Aaron Courville, and Alexandre Lacoste. Mastering the unsupervised reinforcement learning benchmark from pixels. In <i>International Conference on Machine Learning</i> , pp. 28598–28617. PMLR, 2023.
718 719 720	Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In <i>International Conference on Learning Representations</i> , 2020.
721 722 723 724	David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. <i>Science</i> , 362(6419):1140–1144, 2018.
725 726 727	Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information estimators. In <i>International Conference on Learning Representations</i> , 2020.
728 729	Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for markov decision processes. <i>Journal of Computer and System Sciences</i> , 74(8):1309–1331, 2008.
730 731 732 733	DJ Strouse, Kate Baumli, David Warde-Farley, Volodymyr Mnih, and Steven Stenberg Hansen. Learning more skills through optimistic exploration. In <i>International Conference on Learning</i> <i>Representations</i> , 2022.
734 735 736	Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud- den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. <i>arXiv</i> <i>preprint arXiv:1801.00690</i> , 2018.
737 738 739 740	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
741 742 743	Jingda Wu, Zhiyu Huang, and Chen Lv. Uncertainty-aware model-based reinforcement learning: Methodology and application in autonomous driving. <i>IEEE Transactions on Intelligent Vehicles</i> , 8(1):194–203, 2022.
744 745 746	Ruihan Yang, Huazhe Xu, Yi Wu, and Xiaolong Wang. Multi-task reinforcement learning with soft modularization. <i>Advances in Neural Information Processing Systems</i> , 33:4767–4777, 2020.
747 748 749	Rushuai Yang, Chenjia Bai, Hongyi Guo, Siyuan Li, Bin Zhao, Zhen Wang, Peng Liu, and Xuelong Li. Behavior contrastive learning for unsupervised skill discovery. In <i>Proceedings of the 40th</i> <i>International Conference on Machine Learning</i> , pp. 39183–39204, 2023.
750 751 752	Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with pro- totypical representations. In <i>International Conference on Machine Learning</i> , pp. 11920–11931. PMLR, 2021.
753 754 755	Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous con- trol: Improved data-augmented reinforcement learning. In <i>International Conference on Learning</i> <i>Representations</i> , 2022. URL https://openreview.net/forum?id=_SJyyes8.

Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data. *Advances in Neural Information Processing Systems*, 34:25476–25488, 2021.

Yifu Yuan, Jianye HAO, Fei Ni, Yao Mu, YAN ZHENG, Yujing Hu, Jinyi Liu, Yingfeng Chen, and Changjie Fan. EUCLID: Towards efficient unsupervised reinforcement learning with multichoice dynamics model. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=xQAjSr64PTc.

Weitong ZHANG, Dongruo Zhou, and Quanquan Gu. Reward-free model-based reinforcement learning with linear function approximation. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, 2021. URL https://openreview.net/forum?id=IoEnnwAP7aP.

A THEORETICAL PROOFS

A.1 PROOF OF EQ. 2

Proof. As we discussed in 3.1, since we uniformly sample skills from the skill set that contains n skills, we have p(z) = 1/n for each skill. Then we have

$$I_{\rm SD3}(s,z) = \log \frac{\lambda n d_z^{\pi}(s)}{\lambda d_z^{\pi}(s) + \sum_{z' \neq z} d_{z'}^{\pi}(s)} = \log \frac{\lambda n d_z^{\pi}(s)}{\lambda d_z^{\pi}(s) + \rho_{z^c}(s)}.$$
 (10)

Then, the gradient of $I_{SD3}(s, z)$ to $\rho_{z^c}(s)$ becomes

$$\nabla_{\rho_{z^{c}}} I_{\text{SD3}}(s, z) = \frac{\lambda d_{z}^{\pi}(s) + \rho_{z^{c}}(s)}{\lambda n d_{z}^{\pi}(s)} \frac{-\lambda n d_{z}^{\pi}(s)}{\left(\lambda d_{z}^{\pi}(s) + \rho_{z^{c}}(s)\right)^{2}} = -\frac{1}{\lambda d_{z}^{\pi}(s) + \rho_{z^{c}}(s)}.$$
(11)

This completes the proof.

A.2 PROOF OF THEOREM 3.1

Proof. For clarity, we write I_{SD3} as $I_{SD3}(\lambda)$ to explicitly highlight its dependency on the parameter λ in the following context. We first note that the function $I_{SD3}(\lambda)$ is monotonically increasing relative to λ , and $I_{SD3}(\lambda)$ is equal to I(S; Z) when $\lambda = 1$. Therefore, the first inequality

$$I(S;Z) \leq I_{\text{SD3}}(\lambda)$$

always holds for $\lambda \ge 1$. Next, it remains to prove the second inequality, which suffices to give an upper bound of $I_{SD3}(\lambda) - I(S; Z)$. Note that

$$I_{\mathrm{SD3}}(\lambda) - I(S;Z)$$

$$= \mathbb{E}_{z \sim p(z), s \sim d_{z}^{\pi}(s)} \left[\log \frac{\lambda \, d_{z}^{\pi}(s)}{\lambda \, d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')} \cdot \frac{d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')}{d_{z}^{\pi}(s)} \right]$$

$$= \log \lambda + \mathbb{E}_{z \sim p(z), s \sim d_{z}^{\pi}(s)} \left[\log \frac{d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')}{\lambda \, d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')} \right]$$

$$= \log \lambda - \mathbb{E}_{z \sim p(z), s \sim d_{z}^{\pi}(s)} \left[\log \frac{d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z') + (\lambda - 1)d_{z}^{\pi}(s)p(z)}{d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')} \right]$$

$$= \log \lambda - \mathbb{E}_{z \sim p(z), s \sim d_{z}^{\pi}(s)} \left[\log \left(1 + (\lambda - 1) \frac{d_{z}^{\pi}(s)p(z)}{d_{z}^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')} \right) \right].$$
(12)

Recalling that $d_z^{\pi}(\cdot)$ denotes the state density of skill z, and p(z) is the probability density function of skill z, we know that the term

808
809
$$\log\left(1 + (\lambda - 1)\frac{d_z^{\pi}(s)p(z)}{d_z^{\pi}(s)p(z) + \sum_{z' \neq z} d_{z'}^{\pi}(s)p(z')}\right)$$
(13)

is always non-negative for $\lambda \geq 1$. Therefore, we have

$$I_{SD3}(\lambda) \le \log \lambda + I(S; Z). \tag{14}$$

This completes the proof.

Proof. In this proof, we first give a formulation of the intrinsic reward in a linear parameterized assumption, and then discuss the special case of tabular MDPs.

With linear assumptions, we denote $\eta(s_t, z_t) \in \mathbb{R}^d$ as the feature vector of (s_t, z_t) , which is ex-tracted by the encoder network of CVAE. The decoder network is assumed to be a linear function of the feature vector as $\hat{s}_t = W_t \eta(s_t, z_t)$, where $W_t \in \mathbb{R}^{c \times d}$ and $\hat{s}_t \in \mathbb{R}^c$. Then the reconstruction of the state becomes a regularized least-squared problem that captures the prediction error given a dataset \mathcal{D}_m , where m is the number of episodes in the dataset. Thus, we have

$$W_t = \arg\min_{W} \sum_{i=0}^{m} \left\| s_t^i - W\eta(s_t^i, z_t^i) \right\|_F^2 + \kappa \cdot \|W\|_F^2,$$
(15)

where $\|\cdot\|_F$ denotes the Frobenius norm. We further define the following noise with respect to the least-square problem in Eq. (15) as

$$s_t = W_t \eta(s_t, z_t) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}).$$
 (16)

Here we consider the estimation error ϵ in Eq. (15) to follow the standard multivariate Gaussian distribution.

Recall that our practical intra-skill exploration reward is $D_{\rm KL}[Q_{\phi}(h|s,z)|| r(h)]$, where Q_{ϕ} is a posterior network compressing the representation of each state and skill with parameter ϕ , and r(h)is the marginal distribution of the latent variable, where we follow previous works (Alemi et al., 2017; Bai et al., 2021) to consider the marginal as the standard normal distribution. Then we re-define the intrinsic reward in a Bayesian perspective, where we introduce Φ to denote the total parameters, as

$$r_z^{\exp}(s) = \mathbb{E}_{\Phi} D_{\mathrm{KL}}[Q_{\phi}(h|s,z) \parallel r(h)] = \mathcal{H}(Q^{\mathrm{margin}}) - \mathcal{H}(Q_{\phi}(h|s,z)), \tag{17}$$

where $Q^{\text{margin}} = Q(s, z) | \mathcal{D}_m$ is the margin distribution of the encoding over the posterior of the parameters Φ . In practice, we replace the expectation over posterior Φ by the corresponding point estimation, namely the parameter ϕ of the neural networks trained with SD3 model on the dataset \mathcal{D}_m . Formally, considering the Bayesian form of learning objective, we have

$$r_z^{\exp}(s) = \mathcal{H}(Q^{\mathrm{margin}}) - \mathcal{H}(Q_\phi(h|s,z)) = \mathcal{H}(Q(s,z,S)|\mathcal{D}_m) - \mathcal{H}(Q(s,z,S)|\Phi,\mathcal{D}_m), \quad (18)$$

where Q is a neural network in practice. We adopt the mapping

$$Q(s, z, S)|\Phi, \mathcal{D}_m = Q_\phi(h|s, z) \tag{19}$$

since Q_{ϕ} is trained to reconstruct the variable S, where ϕ constitutes a part of the parameters of the total parameters Φ . According to Data Processing Inequality, the post-processing of the signal does not increase information, and we can understand Q as post-processing mapping the state-skill vector via an encoder network. Then we have the following inequality for the information-gain term:

$$r_{z}^{\exp}(s) = \mathcal{H}(Q(s,z,S)|\mathcal{D}_{m}) - \mathcal{H}(Q(s,z,S)|\Phi,\mathcal{D}_{m}) \\ \leq \mathcal{H}(s,z,S|\mathcal{D}_{m}) - \mathcal{H}(s,z,S|\Phi,\mathcal{D}_{m}) = I(\Phi;(s,z,S)|\mathcal{D}_{m}),$$
(20)

where we denote (s, z) as realizations as they are sampled from the dataset as input, and S is a random variable that is learned to reconstruct by parameter Φ . The inequality can be tight since $Q(\cdot)$ is trained by reconstruction, which contains sufficient information about (s, z, S).

In the following, we will prove the following inequality in a linear case with a parameter W_t con-sidered in Eq. (15), as

$$r_{z_t}^{\exp}(s_t) \le I(W_t; (s_t, z_t, S_t) | \mathcal{D}_m) \le \frac{c}{2} [\eta(s_t, z_t)^\top \Lambda_t^{-1} \eta(s_t, z_t)],$$
(21)

the $[\eta(s_t, z_t)^{\top} \Lambda_t^{-1} \eta(s_t, z_t)]$ term is known as an upper-confidence-bound (UCB)-term in linear MDPs (Jin et al., 2023; Cai et al., 2020), and $\Lambda_t = \sum_{j=1}^m \eta(s_j, z_j) \eta(s_j, z_j)^{\top} + \kappa \cdot \mathbf{I}$ is the covari-ance matrix of the samples in the dataset. Finally, we will connect the UCB-term to the count-based bonus in the tabular case.

Let denote $vec(W_t)$ as vectorization of $W_t \in \mathbb{R}^{c \times d}$, and also $\tilde{\eta}(s_t, z_t)$

$$\begin{array}{c}
\mathbf{870}\\
\mathbf{871}\\
\mathbf{872}\\
\mathbf{873}\\
\mathbf{873}\\
\mathbf{874}\\
\mathbf{875}\\
\mathbf{876}\\
\mathbf{876}\\
\mathbf{876}\\
\mathbf{876}\\
\mathbf{877}\\
\mathbf{877}\\
\mathbf{878}\\
\mathbf{878}\\
\mathbf{880}\\
\mathbf{881}\\
\mathbf{881}\\
\mathbf{881}\\
\mathbf{876}\\
\mathbf{876}$$

then it is not difficult to verify that $\operatorname{vec}(W_t)^{\top} \tilde{\eta}(s_t, z_t) = W_t \eta(s_t, z_t)$. By the definition of the mutual information, we observe

$$I(W_t; [s_t, z_t, S_t] \mid \mathcal{D}_m) = I(\operatorname{vec}(W_t); [s_t, z_t, S_t] \mid \mathcal{D}_m)$$

= $\mathcal{H}(\operatorname{vec}(W_t) \mid \mathcal{D}_m) - \mathcal{H}(\operatorname{vec}(W_t) \mid \mathcal{D}_m \cup (s_t, z_t, S_t))$
= $\frac{1}{2} \log \det \left(\operatorname{Var}(\operatorname{vec}(W_t) \mid \mathcal{D}_m) \right) - \frac{1}{2} \log \det \left(\operatorname{Var}(\operatorname{vec}(W_t) \mid \mathcal{D}_m \cup (s_t, z_t, S_t)) \right).$ (23)

Next, we need to obtain $\operatorname{Var}(\operatorname{vec}(W_t) \mid \mathcal{D}_m)$ and $\operatorname{Var}(\operatorname{vec}(W_t) \mid \mathcal{D}_m \cup (s_t, z_t, S_t))$. Recalling that ϵ satisfies the standard Gaussian distribution in Eq. (16), we can conclude that

$$s_t | \eta_t, W_t \sim \mathcal{N}(\operatorname{vec}(W_t)^\top \tilde{\eta}(s_t, z_t), \mathbf{I})$$

Assuming the prior distribution $W \sim \mathcal{N}(0, \mathbf{I}/\kappa)$, then the prior of $\operatorname{vec}(W)$ also follows from $\mathcal{N}(0, \mathbf{I}/\kappa)$. Moreover, using Bayes' theorem and plugging the probability of $p(\operatorname{vec}(W_t))$, we have

$$\log p(\operatorname{vec}(W_t) \mid \mathcal{D}_m) = \log p(\operatorname{vec}(W_t)) + \log p(\mathcal{D}_m \mid \operatorname{vec}(W_t)) - \log p(\mathcal{D}_m) = -\|\operatorname{vec}(W_t)\|^2 / 2 - \sum_{i=1}^m \|\operatorname{vec}(W_t)\tilde{\eta}(s_t^i, z_t^i) - s_{t+1}^i\|^2 / 2 + \operatorname{Const}$$
(24)
$$= -(\operatorname{vec}(W_t) - \tilde{\mu}_{t,m})^\top \tilde{\Lambda}_{t,m}^{-1} (\operatorname{vec}(W_t) - \tilde{\mu}_{t,m}) / 2 + \operatorname{Const},$$

where $\tilde{\mu}_t$ and Λ_t in the last equality are defined as

$$\tilde{\mu}_{t,m} = \tilde{\Lambda}_t^{-1} \sum_{i=0}^m \tilde{\eta}(s_t^i, z_t^i) s_{t+1}^i \in \mathbb{R}^{cd}, \qquad \tilde{\Lambda}_{t,m} = \sum_{i=0}^m \tilde{\eta}(s_t^i, z_t^i) \tilde{\eta}(x_t^i, z_t^i)^\top + \kappa \cdot \mathbf{I} \in \mathbb{R}^{cd \times cd}.$$

Taking the left-hand side of log to the right. The Eq. (24) implies the distribution of $vec(W_t)$ $\mathcal{D}_m \sim N(\tilde{\mu}_{t,m}, \tilde{\Lambda}_{t,m}^{-1})$. Hence, we can get

$$\operatorname{Var}(\operatorname{vec}(W_t) \mid \mathcal{D}_m) = \tilde{\Lambda}_{t,m}^{-1}, \qquad \operatorname{Var}(\operatorname{vec}(W_t) \mid \mathcal{D}_m \cup (s_t, z_t, S_t)) = \tilde{\Lambda}_{t,m+1}^{-1}.$$
(25)

We proceed to derive Eq. (23) by applying Eq. (25), from which we obtain

$$I(\operatorname{vec}(W_t); [s_t, z_t, S_{t+1}] | \mathcal{D}_m) = \frac{1}{2} \log \det \left(\tilde{\Lambda}_{t,m}^{-1} \right) - \frac{1}{2} \log \det \left(\tilde{\Lambda}_{t,m+1}^{-1} \right)$$
$$= \frac{1}{2} \log \det \left(\tilde{\Lambda}_{t,m+1} + \tilde{\eta}(s_t, z_t) \tilde{\eta}(s_t, z_t)^\top \right) - \frac{1}{2} \log \det \left(\tilde{\Lambda}_{t,m} \right)$$
$$= \frac{1}{2} \log \det \left(\tilde{\eta}(s_t, z_t)^\top \tilde{\Lambda}_t^{-1} \tilde{\eta}(s_t, z_t) + \mathbf{I} \right),$$
(26)

where the last equality holds by applying the Matrix Determinant Lemma to the first term. Recalling our definition of $\tilde{\eta}(s_t, z_t)$, the state-skill pairs are finite in the tabular case, so we have

$$\tilde{\Lambda}_{t} = \sum_{i=0}^{m} \tilde{\eta}(s_{t}^{i}, z_{t}^{i}) \tilde{\eta}(x_{t}^{i}, z_{t}^{i})^{\top} + \kappa \cdot \mathbf{I} \\
= \begin{bmatrix} \sum \eta(s_{0}, z_{0}) \eta(s_{0}, z_{0})^{\top} + \kappa I & 0 & \cdots & 0 \\ 0 & \sum \eta(s_{1}, z_{1}) \eta(s_{1}, z_{1})^{\top} + \kappa I & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sum \eta(s_{m}, z_{m}) \eta(s_{m}, z_{m})^{\top} + \kappa I \end{bmatrix}.$$
(27)

Then, $\tilde{\eta}(s_t, z_t)^{\top} \tilde{\Lambda}_t^{-1} \tilde{\eta}(s_t, z_t)$ can be rewritten as

m

$$\begin{split} \tilde{\eta}(s_{t}, z_{t})^{\top} \Lambda_{t}^{-1} \tilde{\eta}(s_{t}, z_{t}) \\ &= \begin{bmatrix} \eta(s_{t}, z_{t})^{\top} & 0 & \cdots & 0 \\ 0 & \eta(s_{t}, z_{t})^{\top} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \eta(s_{t}, z_{t})^{\top} \end{bmatrix} \begin{bmatrix} \Lambda^{-1} & 0 & \cdots & 0 \\ 0 & \Lambda^{-1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \Lambda^{-1} \end{bmatrix} \begin{bmatrix} \eta(s_{t}, z_{t}) & 0 & \cdots & 0 \\ 0 & \eta(s_{t}, z_{t}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \eta(s_{t}, z_{t}) \end{bmatrix}$$
(28)
$$= \begin{bmatrix} \eta(s_{t}, z_{t})^{\top} \Lambda^{-1} \eta(s_{t}, z_{t}) & 0 & \cdots & 0 \\ 0 & \eta(s_{t}, z_{t})^{\top} \Lambda^{-1} \eta(s_{t}, z_{t}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \eta(s_{t}, z_{t})^{\top} \Lambda^{-1} \eta(s_{t}, z_{t}) \end{bmatrix} \in \mathbb{R}^{c \times c}. \end{split}$$

Therefore, by eliminating the determinant based on the expression in Eq. (28) and applying the inequality $\log(1+x) \le x$ for $x \ge 0$, we can further bound Eq (26) from above as

$$I(\operatorname{vec}(W_t); [s_t, z_t, S_{t+1}] | \mathcal{D}_m) = \frac{1}{2} \cdot \log \det \left(\tilde{\eta}(s_t, z_t)^\top \tilde{\Lambda}_t^{-1} \tilde{\eta}(s_t, z_t) + \mathbf{I} \right)$$
$$= \frac{c}{2} \cdot \log \left(\eta(s_t, z_t)^\top \Lambda^{-1} \eta(s_t, z_t) + 1 \right)$$
$$\leq \frac{c}{2} \cdot \eta(s_t, z_t)^\top \Lambda^{-1} \eta(s_t, z_t).$$
(29)

Hence, based on Eq. (20) and Eq. (29), we conclude that

$$r_{z}^{\exp}(s_{t}) \leq I(W_{t}; [s_{t}, z_{t}, S_{t+1}] | \mathcal{D}_{m}) = I(\operatorname{vec}(W_{t}); [s_{t}, z_{t}, S_{t+1}] | \mathcal{D}_{m}) \leq \frac{c}{2} \cdot \eta(s_{t}, z_{t})^{\top} \Lambda^{-1} \eta(s_{t}, z_{t}).$$
(30)

In tabular cases (Auer & Ortner, 2006), the state and skill are considered as finite and countable. Let $d = |S| \times |Z|$. Recall that $\eta(s_t, z_t) \in \mathcal{R}^{|S||Z|}$ is the one-hot vector with a value of 1 at position $(s_t, z_t) \in S \times Z$, i.e.,

$$\eta(s_j, z_j) = \begin{bmatrix} 0\\ \vdots\\ 1\\ \vdots\\ 0 \end{bmatrix} \in \mathbb{R}^d, \text{ and } \eta(s_j, z_j)\eta(s_j, z_j)^\top = \begin{bmatrix} 0 & \cdots & 0 & \cdots & 0\\ \vdots & \ddots & & \vdots\\ 0 & & 1 & & 0\\ \vdots & & \ddots & \vdots\\ 0 & \cdots & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{d \times d}.$$
(31)

We denote the gram matrix $\Lambda_j = \sum_{i=0}^m \eta(s_j^i, z_j^i) \eta(s_j^i, z_j^i)^\top + \kappa \cdot \mathbf{I}$ for $\kappa > 0$ as covariance matrix given a dataset \mathcal{D}_m . Since we denote η as a one-hot vector, and Λ as the sum of all the matrices $\eta(s_j, z_j) \eta(s_j, z_j)^\top$, each diagonal element of Λ can be seen as the corresponding count $N(s_j, z_j)$ for the state-skill pair, i.e.

966	$\Gamma N(s_0, z_0) + \kappa$	0		0 J
967	0	$N(s_1, z_1) + \kappa$		0
968	:	·		:
969 Λ =	$= \begin{bmatrix} & \cdot \\ & 0 \end{bmatrix}$	N($(s_i, z_i) + \kappa$	0
970	· ·	`	•	
971			·. 	$N(s_{m}, z_{m}) + \kappa$

Moreover, given a dataset, the expression on the right side of the theorem's inequality is inversely
 proportional to the total number of state-skill pairs; in other words,

$$\eta(s_j, z_j)^{\top} \Lambda_t^{-1} \eta(s_j, z_j) = \frac{1}{N(s_j, z_j) + \kappa}.$$
(32)

According to Eq. (21), we have the following relationship in the tabular case:

$$r_{z_t}^{\exp}(s_t) \le I(W_t; (s_t, z_t, S_t) | \mathcal{D}_m) \le \frac{c}{2} [\eta(s_t, z_t)^\top \Lambda_t^{-1} \eta(s_t, z_t)] = \frac{|\mathcal{S}|/2}{N(s_t, z_t) + \kappa}.$$
 (33)

The first inequality is due to the Data Processing Inequality according to Eq. (20). The bound is tight since $Q(\cdot)$ is trained by reconstruction, which contains sufficient information about (s, z, S). The second inequality is tight when $\eta(s_t, z_t)^{\top} \Lambda_t^{-1} \eta(s_t, z_t) \to 0$, which means that the count of state-action pair is large. In the last equation, c is the count of all states in the tabular space. Thus, we have

$$r_z^{\exp}(s) \approx \frac{|\mathcal{S}|/2}{N(s,z) + \kappa},\tag{34}$$

if the count of N(s, z) is large. Intuitively, optimizing the reward $\eta(s, z)^{\top} \Lambda^{-1} \eta(s, z)$ incentivizes the agent to increase the visitation of (s, z). Furthermore, since we have proven that Eq. (21) holds, we can state that in the tabular case, maximizing the intra-skill reward is equivalent to maximizing the count-based rewards (Bellemare et al., 2016; Ostrovski et al., 2017). The intra-skill exploration reward encourages the skill-conditional policy to increase the visitation times of those rare state-skill pairs.

B Hyper-parameters and Implementation Details

999 B.1 HYPER-PARAMETERS

1000 We utilize the baselines from the open-source implementations of URLB (https: 1001 //github.com/rll-research/url_benchmark), CIC (https://github.com/ 1002 rll-research/cic), and BeCL (https://github.com/Rooshy-yang/BeCL), 1003 keeping their hyper-parameters fixed throughout both the pre-training and fine-tuning 1004 stages. For CSD and Metra, due to the absence of their experiments in state- and pixel-1005 based URLBs, we re-implement them in these benchmarks on their official implemen-(CSD https://github.com/seohongpark/CSD-locomotion, tations Metra https://github.com/seohongpark/METRA). Table 1 details the hyper-parameters used for SD3 and DDPG. 1008

1009

1011

1020 1021

1010 B.2 IMPLEMENTATION DETAILS

Soft Modularized CVAE To achieve SD3, we utilize a soft modularized CVAE to estimate the 1012 state density $d_z^{\pi}(s)$ of one skill. Specifically, we forward the state s through an MLP or CNN to 1013 obtain a d-dimensional state embedding $f_1(s)$ and, similarly, obtain a d-dimensional skill embedding 1014 $f_2(z)$. We use $f_1(s)$ as the input to the unconditional basic network, and $f_1(s) \odot f_2(z)$ as the input 1015 to the routing network. The basic network comprises n layers, each containing m modules, for 1016 progressively extracting features. The routing network contains n-1 gating layers, which provide 1017 a probability vector p^l based on the input as shown in Eq. (4) to weight the contribution of the *l*-th 1018 layer's modules to the l + 1-th layer's modules. Particularly, the probability vector which outputs 1019 from the first layer of the routing network is represented as

$$p^{l=1} = \mathcal{W}^l \big(\text{ReLU}(f_1(s) \odot f_2(z)) \big).$$
(35)

Then, the probability vector is normalized using the softmax function as \hat{p}^l and the input to each module in the basic network can be expressed as

1025
$$g_i^{l+1} = \sum_j \hat{p}_{i,j}^l (\text{ReLU}(\mathcal{W}_j^l g_j^l)), \qquad (36)$$

974 975

976 977

987

988

989

990

991

992

993

994

995 996

997

998

027		
028	SD3 hyper-parameter	Value
)29	Skill dim	16 discrete
130	Softmax Temperature T	1
0.4	Skill sampling frequency (steps)	50
J31	Exploration ratio α	$\{0.04, 2.0\}$
)32	Weight Parameter λ	1.5
)33	CVAE Encoder arch.	$\dim(S) \rightarrow 1024 \rightarrow 1024 \rightarrow 1024 \rightarrow 40 * 2 \text{ ReLU (MLP)}$
)34	CVAE Decoder arch.	$40 \rightarrow 1024 \rightarrow 1024 \rightarrow 1024 \rightarrow \dim(S) \text{ ReLU (MLP)}$
)35	DDPG hyper-parameter	Value
126	Replay buffer capacity	10^{6}
50	Action repeat	1
)37	Seed frames	4000
)38	<i>n</i> -step returns	3
)39	Mini-batch size	1024
040	Seed frames	4000
)/1-1	Discount (γ)	0.99
/41	Optimizer	Adam
J42	Learning rate	10^{-4}
)43	Agent update frequency	2
)44	Critic target EMA rate (τ_Q)	0.01
)45	Features dim.	1024
146	Hidden dim.	1024
	Exploration stddev clip	0.3
J4/	Exploration stddev value	0.2
048	Number pretraining frames	$2 imes 10^6$
049	Number finetuning frames	1×10^5
050		

1026

1007

where $\hat{p}_{i,j}^l$ weights the *j*-th module in the *l*-th layer to contribute to the *i*-th module in the *l* + 1-th layer, g_j^l is the input to the *j*-th module in the *l*-th layer and \mathcal{W}_j^l represents the module parameters.

1055 By progressively extracting features of state s while incorporating the weight information, the en-1056 coder transforms the state s into the mean $\mu(s|z)$ and variance $\sigma^2(s|z)$ of the latent space condi-1057 tioned on the skill z. The latent representation h is generated using the reparameterization trick, ensuring gradients can be backpropagated through the sampling process. Specifically, this is done 1058 as $h = \mu + \sigma \cdot \epsilon$, where ϵ is the noise sampled from a standard Gaussian distribution. The decoder 1059 then progressively up-samples and reconstructs the output state \hat{s} from the latent representation h. 1060 incorporating the weight information generated by the routing network. We train the entire soft 1061 modularized CVAE by maximizing \mathcal{L}^{elbo} as given in Eq. (3), enabling it to more accurately estimate 1062 $d_z^{\pi}(s).$ 1063

1003

1064

Practical Implementation We propose the complete SD3 algorithm in Algorithm 2. We conduct our experiments using an RTX 4090 GPU. Each run in the state-based URLB environment takes approximately 1 day, while runs in the Maze environment requires about 3 hours each. For the pixel-based URLB environment, each run takes around 4 days or less.

- 1069 1070
- 1071

C VISUALIZATION

1072 1073

1074 C.1 TREE-LIKE MAZE

1075

As shown in Figure 6, we conduct additional experiments in the tree-like maze to visualize the skills learned by SD3. It can be observed that DIAYN and DADS only reach the middle of the maze, whereas SD3 successfully reaches the bottom of the maze. The proposed latent space reward in SD3 demonstrates strong exploration ability in large-scale mazes. Moreover, the trajectories of different skills remain distinguishable in SD3.

Alg	prithm 2: Complete SD3 algorithm
Iı	put: number of pre-training frames N_{PT} , number of fine-tuning frames N_{FT} , batch size N,
sł	ill sampling frequency N_{update} , skill set \mathcal{Z} , exploration ratio α .
h	itialize Environment, CVAE Q_{ϕ} , actor π_{θ} , critic Q_{φ} , replay buffer \mathcal{D} .
//. £.	Pre-training
10	$\mathbf{r} t = 1$ to N_{PT} do Randomly choose a from the skill set \mathcal{T} every N_{res} steps
	Interact with environment by $\pi_0(a e z)$
	Store the transition in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \sqcup (s_{+}, a_{+}, s'_{-}, z)$
	if $t > 4,000$ then
	Sample a batch from $\mathcal{D}: \{s_t, a_t, s'_t, z\}^N \sim \mathcal{D}.$
	Update CVAE Q_{ϕ} via $\mathcal{L}^{\text{elbo}}$ in Eq. (3).
	Use CVAE Q_{ϕ} to compute $d_z^{\pi}(s)$ and $d_{z'\neq z}^{\pi}(s)$.
	Compute $r_z^{sd3}(s)$ and $r_z^{exp}(s)$ with Eqs. (6)-(7).
	Compute the intrinsic reward $r^{\text{int}} = r_z^{\text{sd3}}(s) + \alpha \cdot r_z^{\text{exp}}(s)$
	Update actor π_{θ} and critic Q_{φ} using intrinsic reward r^{int} .
	end if
e	nd for
//. £c	Fine-tuning $r t = 1$ to N do
ю	$T t = 1$ to N_{FT} to Use pre-training models to initialize actor π_{ee} and critic O_{ee}
	Randomly sample a skill z^* from Z and fix the z^*
	Interact with environment by $\pi_{a'}$.
	Store the transition in replay buffer $\mathcal{D} \leftarrow \mathcal{D} \cup (s_t, a_t, r_t, s'_t, z^*)$.
	if $t \ge 4,000$ then
	Sample a batch from $\mathcal{D}: \{s_t, a_t, r_t, s'_t, z^*\}^N \sim \mathcal{D}.$
	Use extrinsic reward r_t obtained from downstream task to update $\pi_{\theta'}$ and $Q_{\varphi'}$.
	end if
e	nd for
C^{2}	DEEDMIND CONTROL SUITE
C.2	DEEL MIND CONTROL SUITE
Figu	re 7 shows the learned skills in the Walker, Quadruped, and Jaco Arm domains. The result
show	vs SD3 can learn various locomotion skills, including standing, walking, rolling, moving, and
som	ersault; and also learns various manipulation skills by moving the arm to explore different areas,
oper	ning and closing the gripper in different locations. The learned meaningful skills lead to superior
gene	ralization performance in the fine-tuning stage of various downstream tasks.
D	ABLATION STUDIES
ν	ADEALION STODIES
D 1	ΤΗΕ ΕΧΡΙ ΟΡΑΤΙΟΝ ΒΑΤΙΟ
D.1	
We	conduct an ablation on the different exploration ratios α , Specifically, with the hyper-parameter
α , the	he reward is represented as:
	$r^{\text{total}}(e) = r^{\text{sd3}}(e) + \alpha \cdot r^{\exp}(e) $ (37)
	$T_{z} = (s) - T_{z} = (s) + \alpha \cdot T_{z} = (s).$ (37)
As i	llustrated in Figure 8(a), when α is set to 0 and 0.02, the agent can learn distinguishable and
conv	regent skills but fails to fully explore the maze. When α is set to 0.08, the agent explores
suffi	ciently, but the trajectories at the endpoints are quite scattered, indicating that the learned skill
	ciently, but the trajectories at the endpoints are quite scattered, indicating that the feather skin
strat	egies lack stability. In contrast, $\alpha = 0.04$ balances exploration and the skill diversity.
strat	egies lack stability. In contrast, $\alpha = 0.04$ balances exploration and the skill diversity.

1131 According to our analysis, when the proportion of exploration is deficient or even absent, SD3 solely 1132 maximize I_{SD3} . Conversely, an excessively high α can overly prioritize intra-skill exploration, 1133 resulting in instability within the learned skills. Empirically, we have found that $\alpha = 0.04$ can lead to promising results in downstream tasks in the Quadruped domain.



Figure 6: Additional experiments in the tree-like Maze with different numbers of skills. Under different environmental conditions, SD3 demonstrates superior exploration capabilities while still learning distinguishable skills, outperforming DIAYN and DADS.



Figure 7: Skill visualization in DMC. It can be observed that SD3 learns dynamic and valuable skills, which enable the agent to quickly adapt to downstream tasks.

- 1185 1186
- 1187



Figure 8: Results for the impact of exploration ratio. (a) We conduct experiments with different α in the maze and found that varying α values significantly impact both the state coverage and the stability of learned skills. (b) In the *Quadruped* domain, different α also have a notable effect on the performance of various downstream tasks.



Figure 9: Ablation on the soft modularization structure.

1223 D.2 IMPACT OF SOFT MODULARIZATION

As mentioned in section 3.1, we use CVAE to estimate the state density of different skills. To enhance the accuracy of estimation in complex state spaces, we have introduced soft modularization into the traditional CVAE structure. Consequently, we conduct an ablation study on the soft modularization. Aggregated scores are reported in Figure 9. We observe that SD3 with soft modularized CVAE obtains superior performance, as it has sufficient capacity to learn the density information of different skills for the same state in complex state spaces, while the skill density estimation of one skill may intervene with those of other skills in the traditional CVAE.

D.3 TEMPERATURE IN SOFTMAX

In section 3.1, we mentioned that the normalized weight $\hat{p}_{i,j}^l$ for the routing network is computed with the equation $\hat{p}_{i,j}^l = \exp(p_{i,j}^l)/(\sum_{j=1}^m \exp(p_{i,j}^l))$. This is implemented using *Softmax*, where we follow the previous work (Hinton et al., 2015) to introduce a *temperature* T to control the level of uncertainty in output probabilities. The formula is as follows:

$$\hat{p}_{i,j}^{l} = \frac{\exp(p_{i,j}^{l}/T)}{\sum_{j=1}^{m} \exp(p_{i,j}^{l}/T)}.$$
(38)

From the above formula, it can be observed that when the temperature T = 1, it resembles the original softmax function. As T decreases, the distribution output by softmax gradually becomes



Figure 10: Results for the impact of softmax temperature. (a) We exhibit the performance of the agent with different temperatures in the *Quadruped* domain. (b) When the temperature is set to 100, SD3 becomes a count-based pure exploration method. It demonstrates a certain degree of environment exploration capability but lacks empowerment in the environment.

more extreme, eventually converging to a deterministic distribution. Conversely, as *T* increases, the softmax gradually tends to derive a uniform distribution. Here we perform the ablation study on the temperature coefficient.

The result is illustrated in Figure 10(a). When T values are 0.1 and 0.01, the output of the softmax function in the routing network will gradually approximate $\operatorname{argmax}(p_i^l)$, at which point each training iteration utilizes only a single module from each layer. This practice inevitably diminishes the accuracy of estimating $d_z^{\pi}(s)$; when T values are 10 and 100, the routing network tends to output uniformly distributed weight values, causing the routing network to fail. The entire network structure can be approximated as a VAE composed of multiple modules. Given the loss of skill z information in the basic network, the intrinsic reward of SD3 can be repsented as follows:

$$r_z^{\text{total}}(s) = r_z^{\text{sd3}}(s) + r_z^{\text{exp}}(s)$$

1272

1278

1279

1283 1284 1285 $r_{z}^{\text{cond}}(s) = r_{z}^{\text{cond}}(s) + r_{z}^{\text{cond}}(s)$ $= \log \frac{\lambda d^{\pi}(s)}{\lambda d^{\pi}(s)p(z) + \sum_{z' \neq z} d^{\pi}(s)p(z')} + D_{\text{KL}}[Q_{\phi}(\cdot|s)||r(h)]$ $= \log \frac{\lambda \cdot n}{\lambda + n - 1} + D_{\text{KL}}[Q_{\phi}(\cdot|s)||r(h)]$ $= c + D_{\text{KL}}[Q_{\phi}(\cdot|s)||r(h)],$ (39)

where *c* represents a constant. Furthermore, based on Theorem 3.2, the right-hand side of the equation can be approximated as $\frac{|S|/2}{N(s)+\kappa}$. Substituting into Eq. (39), we can obtain:

$$r_z^{\text{total}}(s) \approx c + \frac{|\mathcal{S}|/2}{N(s) + \kappa}.$$
 (40)

At this point, the SD3 method transforms into a count-based exploration approach. As presented in Figure 10(b), the agent learns extremely dynamic behavior, thereby preventing it from adequately adapting downstream tasks.

1289

1291

1290 D.4 Impact of Weight Parameter λ

1292 The discussion in section 3.1 introduces a weight parameter λ in Eq.(1). To investigate the impact 1293 of λ , we conduct an ablation study by varying λ from [0.5, 1.0, 1.5, 2.0, 3.0]. The results, exhibited 1294 in Figure 11, indicate that the performance of SD3 fluctuates within a narrow range when lambda is 1295 greater than 1. Therefore, we conclude that λ is generally applicable in a wider range, and SD3 is 1296 not sensitive to the parameter when $\lambda \ge 1.5$.

1317





Figure 11: Results for the impact of weight pa-1308 rameter in the *Quadruped*. When λ is set to 0.5 1309 or 1, SD3 performs poorly. However, it is ob-1310 served that increasing lambda beyond 1 does not 1311 significantly impact the performance of SD3. 1312

Figure 12: Skill adaption strategies ablation. We test several adaptation methods in the finetuning phase and find that randomly selecting skills perform comparably to using regressmeta, but employing the meta-controller results in a decline in the performance.

1316 D.5 **SKILL ADAPTION STRATEGIES IN FINE-TUNING**

1318 Previous work (Laskin et al., 2021) has shown that during the fine-tuning phase, performance across 1319 different skills does not always level equally; some skills demonstrate weaker adaptability in down-1320 stream tasks, while others show the opposite. Therefore, we investigate various skill adaptation 1321 methods in the state-based environment to assess their impact on algorithm performance in downstream tasks. 1322

1323 In the experiment described in section 5.2, for a fair comparison, we adhere to the standards set in 1324 the URLB, employing a random sampling skills method during the fine-tuning stage to evaluate the 1325 average performance of skills. Therefore, here we introduce two additional skill adaptation methods: 1326 regress-meta and meta-controller. Regress-meta computes the expected reward for each skill during 1327 the first 4K steps of the fine-tuning phase to determine its skill-value, and then selects the skill with the highest skill-value to perform the downstream task. Meta-controller trains an upper-level 1328 controller $\mu(z|s)$ in the fine-tuning phase to select the most appropriate skill for the current state s, 1329 thereby combining it with the policy $\pi(a|s,z)$ trained in the pre-training phase and optimizing the 1330 high-level policy based on $\pi(a|s) = \sum_{z \in \mathcal{Z}} \mu(z|s) \pi(a|s, z)$. 1331

1332 Results are shown in Figure 12. The performance of using regress-meta to select skills shows im-1333 provements compared to randomly selecting skills in *Quadruped Stand*, Walk, and Run but a slight drop in Quadruped Jump. We attribute this to the fact that regress-meta consistently selects the 1334 skill with the highest expected reward during the initial steps of the fine-tuning phase. While this 1335 approach does increase the probability of choosing a skill with good adaptability, there is also a risk 1336 of choosing a skill that performs well during the initial 4K steps but exhibits mediocre performance 1337 thereafter. In contrast, the meta-controller exhibits relatively poor performance. We hypothesize 1338 that the meta-controller usually requires a large number of examples to train, which is difficult to 1339 converge within the 100K fine-tuning steps. 1340

- 1341
- 1342

1343 E NUMERICAL RESULT

1344 1345

In Table 2 and Table 3, we present the mean normalized scores and standard errors of all algorithms across 12 downstream tasks within the state-based URLB experiments. SD3 demonstrates superior 1347 performance across multiple downstream tasks. In Table 4, we present the results of the pixel-based 1348 URLB experiments. Across 8 downstream tasks, SD3 displays notable competitiveness compared 1349 to other baselines. Additionally, we showcase the results of robustness experiments in Table 5.

Domain	Task	DDPG	CSD	Metra	CIC	BeCL	SD3
	Flip	538±27	615 ± 17	600 ± 48	641 ± 26	611±18	595 ± 25
Wallton	Run	325±25	445±13	302 ± 23	450 ± 19	387 ± 22	451±23
walkel	Stand	899±23	962±7	951 ± 7	959 ± 2	952 ± 2	930±5
	Walk	748±47	857±51	756 ± 67	$903{\pm}21$	883 ± 34	914±11
	Jump	236±48	357±39	300 ± 9	565 ± 44	727 ± 15	676±29
Quadmined	Run	157±31	362 ± 60	276 ± 20	445 ± 36	535 ± 13	471 ± 13
Quadruped	Stand	392±73	455 ± 36	637 ± 85	700 ± 55	875±33	847±17
	Walk	229±57	224 ± 18 200 ± 27 621	621 ± 69	743 ± 68	$752{\pm}40$	
	Reach bottom left	72±22	99±7	143 ± 9	154 ± 6	148 ± 13	151±7
Isso	Reach bottom right	117±18	106 ± 6	142 ± 8	149 ± 4	139 ± 14	152±9
Jaco	Reach top left	116±22	101 ± 7	130 ± 13	$149{\pm}10$	125 ± 10	142 ± 7
	Reach top right	94±18	154 ± 11	158 ± 16	163 ± 9	126 ± 10	152 ± 7

Table 2: Results of SD3 and novel competence-based methods on state-based URLB.

Table 3: Results of other baselines on state-based URLB.

1364										
1005	Domain	Task	ICM	Disagreement	RND	APT	ProtoRL	SMM	DIAYN	APS
1365		Flip	390±10	332±7	506±29	606±30	549±21	500 ± 28	361±10	448±36
1366	Walker	Run	267±23	243±14	403±16	384±31	370 ± 22	395 ± 18	184 ± 23	176 ± 18
1007		Stand	836±34	760 ± 24	901±19	921±15	896 ± 20	886±18	789 ± 48	702 ± 67
1307		Walk	696 ± 46	606 ± 51	783 ± 35	784±52	836 ± 25	792 ± 42	450 ± 37	547 ± 38
1368	Overdenned	Jump	205±47	510±28	626±23	416±54	573±40	167±30	498 ± 45	389±72
1260		Run	125 ± 32	357±24	439±7	303±30	$324{\pm}26$	142 ± 28	347 ± 47	201 ± 40
1309	Quadruped	Stand	$260{\pm}45$	579±64	839 ± 25	582±67	625 ± 76	266 ± 48	718 ± 81	435 ± 68
1370		Walk	153 ± 42	386 ± 51	517 ± 41	582 ± 67	494 ± 64	154±36	506 ± 66	385 ± 76
1371	-	Reach bottom left	88±14	117±9	102±9	143±12	118±7	45±7	20 ± 5	84±5
1071	Jaco	Reach bottom right	99±8	122 ± 5	110±7	138±15	138 ± 8	60 ± 4	17 ± 5	94 ± 8
1372		Reach top left	80±13	121 ± 14	88±13	137±20	134±7	39 ± 5	12 ± 5	74 ± 10
1373		Reach top right	$106{\pm}14$	128 ± 11	99±5	170±7	140 ± 9	32 ± 4	21 ± 3	83±11

MORE DISCUSSIONS F

F.1 THE UNIQUE FAVORABLE PROPERTIES OF SD3

Previous skill discovery methods, such as CIC, APS, and BeCL, also encourage exploration while discovering diverse skills. However, in comparison, SD3 possesses its own distinctive properties.

First, SD3 introduces a new objective for skill discovery, which is not derived from maximizing MI. The core principle of SD3 is to promote deviation in exploration regions across different skills, thereby facilitating more effective skill discovery. Unlike previous methods focusing on maximizing a lower-bound of MI, SD3 uses a novel CVAE architecture for density estimation to directly estimate the original objective. Further, as shown in Theorem 3.1, a qualitative analysis reveals that the previous MI objective is merely a special case of SD3.

Second, different from APS, CIC, and BeCL, which explicitly or implicitly maximizes state entropy for exploration, SD3 adopts a novel exploration strategy that resembles count-based exploration. In section 5.4, we confirm that such UCB-style reward is more robust than entropy-based reward. Meanwhile, this exploration reward can be estimated as an byproduct in from the learned CVAE, avoiding the additional mechanisms compared to other methods.

F.2 THE PERFORMANCE OF SD3 COMPARED TO CIC

The quantitative results in Figs 4 and 5(a) indicate that SD3 and CIC are comparable. While SD3 slightly outperforms CIC, the improvement may not be statistically significant. In fact, in the ex-perimental section of the main text, our focus is on showcasing SD3's overall performance and advantages.

Regarding the skill discovery objective, we believe that evaluating the fine-tuning performance of skills is somewhat limited. As demonstrated in the maze experiment (see Figure 3), although CIC achieves the best state coverage, it learns very disorganized skills with mixed trajectories. While CIC attains high scores after fine-tuning, it fails to reflect the core objective of skill discovery, which aims to learn diverse and distinguishable skills. SD3, on the other hand, excels in discovering easily distinguishable skills and also demonstrates competitive performance in downstream tasks.

Stand

Walk

785±18

 475 ± 55

	Domain	Task	APT	CSD	Metra	CIC	BeCL	SD3
-		Flip	803±26	681 ± 56	665 ± 32	836±12	539±8	864±27
	Wallran	Run	506±4	451 ± 41	454 ± 29	504 ± 21	456 ± 14	543±22
Walke	warker	Stand	961±5	958 ± 13	968 ± 4	973±2	968 ± 4	982±1
		Walk	880±37	948 ± 5	949 ± 3	953±5	939 ± 1	945±3
		Jump	557 ± 67	580 ± 74	677 ± 27	723±16	340 ± 32	729±16
	Quadmined	Run	396±9	390 ± 21	276 ± 46	439±3	162 ± 4	438±16
	Quauruped	Stand	705 10	854 L 20	799 ± 20	972 ± 12	502 56	021 1 2

 854 ± 20

 530 ± 19

Table 4: Results of SD3 and baselines on pixel-based URLB.

Table 5: Results of robustness experiments.

 788 ± 29

 181 ± 39

873±13

672±15

 583 ± 56

 283 ± 39

921±3

 680 ± 43

	CIC	CIC	Performance	SD3	SD3	Performance
Task	(Noisy)	(Normal)	Ratio	(Noisy)	(Normal)	Ratio
walker_flip	511±6	641 ± 26	79.72%	554 ± 24	595 ± 25	93.11%
walker_run	319 ± 20	450 ± 19	70.89%	330 ± 25	451 ± 23	73.17%
walker_stand	845 ± 12	959 ± 2	88.11%	909 ± 11	930 ± 5	97.74%
walker_walk	784 ± 46	903 ± 21	86.82%	877±27	914 ± 11	95.95%
quad_jump	384±61	565 ± 44	67.96%	560 ± 48	676 ± 29	82.84%
quad_run	276 ± 48	445 ± 36	62.02%	421 ± 47	471 ± 13	89.38%
quad_stand	424 ± 25	700 ± 55	60.57%	746 ± 93	847 ± 17	88.07%
quad_walk	356±99	621 ± 69	57.32%	529 ± 55	752 ± 40	70.34%
Average	-	-	71.68%	-	-	86.33%

1425 1426

1404

1405 1406

1407

1408

1409

1410

1411

1412

1413 1414

1427

1434

1436

1443

1444

Additionally, we conduct experiments to confirm that SD3 is more robust than CIC in noisy environ-1428 ments. Our four experiments in the main text complement each other and collectively provide suffi-1429 cient evidence that SD3 demonstrates superior and more comprehensive performance compared to 1430 other methods, including the ability to discover distinguishable skills (i.e., in maze/URLB domains), superior performance in downstream tasks (i.e., in state/pixel URLB), and scalability to large-scale 1431 problems (i.e., pixel-based domains). Therefore, we believe that SD3 will be favored over CIC and 1432 other methods for a wide range of tasks. 1433

1435 F.3 THE INTEGRATION OF EXPLORATION AND DIVERSITY REWARDS DURING TRAINING

1437 As vividly displayed in Figure 2, to better explain the key idea behind our algorithm and to illustrate the skill discovery process of SD3, we describe the learning process in an iterative manner. However, 1438 in practice, we first obtain a combined intrinsic reward $r^{\text{int}} = r^{\text{sd}3} + \alpha r^{\text{exp}}$ of two objectives, and 1439 then adopt DDPG as a backbone RL algorithm to learn the policy. We adopt such an optimization 1440 approach because using a combined reward r^{int} only requires learning a single Q-function, which 1441 is more computationally efficient than an iterative process that requires learning two Q-functions. 1442

1445					
1445	Task	CIC	BeCL	SD3	SD3(w/o soft-modu)
1440	Quad Stand	700 ± 55	875 ± 33	847 ± 17	752 ± 64
1447	Quad Walk	621 ± 69	743 ± 68	752 ± 40	642 ± 80
1448	Quad Run	445 ± 36	535 ± 13	471 ± 13	422 ± 34
1449	Quad Jump	565 ± 44	727 ± 15	676 ± 29	589 ± 45
1450	Walker Stand	959 ± 2	952 ± 2	930 ± 5	910 ± 16
1451	Walker Walk	903 ± 21	883 ± 34	914 ± 11	870 ± 30
1452	Walker Run	450 ± 19	387 ± 22	451 ± 23	409 ± 55
1453	Walker Flip	641 ± 26	611 ± 18	595 ± 25	523 ± 33
1454	Jaco Top Left	149 ± 10	125 ± 10	142 ± 7	125 ± 5
1455	Jaco Top Right	163 ± 9	126 ± 10	152 ± 7	117 ± 5
1456	Jaco Bottom Left	154 ± 6	148 ± 13	151 ± 7	134 ± 8
1457	Jaco Bottom Right	149 ± 4	139 ± 14	152 ± 9	122 ± 8

Table 6: Results of SD3 without soft-modularized CVAE.

1459					
1460	Task	Quad Stand	Quad Walk	Quad Run	Quad Jump
1461	Module $= 2$	777 ± 26	638 ± 45	377 ± 38	499 ± 32
1462	Module $= 3$	862 ± 25	650 ± 29	456 ± 26	590 ± 20
1463	Module $= 5$	781 ± 31	799 ± 31	390 ± 32	541 ± 24
1464	Module $= 6$	680 ± 27	323 ± 32	261 ± 31	375 ± 34
1465	Module = 4	847 ± 17	752 ± 40	471 ± 13	676 ± 29
1466					

Table 7: Results for different number of modules.

1468 G ADDITIONAL ABLATION STUDY

1470 G.1 STATE-BASED URLB WITHOUT SOFT-MODULARIZED CVAE

We conduct experiments without the soft-modularized CVAE in state-based URLB, using a standard
CVAE where the encoder consists of a 4-layer MLP network. The results, shown in the Table 6,
demonstrate that even without the soft-modularized CVAE, our method still achieves competitive
performance on several downstream tasks.

1477 G.2 The number of modules in soft-modularized CVAE

We conduct ablation experiments on the number of modules in the state-based quadruped environment, and the results are shown in Table 7. From the table, it can be observed that the performance differences are minimal when the number of modules is set to 2, 3, or 5. However, when the number is increased to 6, there is a significant performance drop. We attribute this to the difficulty in effectively training the soft-modularized structure as the number of modules becomes too large. The relatively comprehensive performance is achieved when the number of modules is set to 4.