# Understanding Effectiveness of Learning Behavioral Metrics in Deep Reinforcement Learning

Anonymous authors Paper under double-blind review

Keywords: behavioral metrics, bisimulation metrics, representation learning, evaluation

# Summary

A key approach to state abstraction is approximating behavioral metrics (notably, bisimulation metrics) in the observation space, and embed these learned distances in the representation space. While promising for robustness to task-irrelevant noise shown in prior work, accurately estimating these metrics remains challenging, requiring various design choices that create gaps between theory and practice. Prior evaluations focus mainly on final returns, leaving the quality of learned metrics and the source of performance gains unclear. To systematically assess how metric learning works in deep RL, we evaluate five recent approaches. We unify them under isometric embedding, identify key design choices, and benchmark them with baselines across 20 state-based and 14 pixel-based tasks, spanning 250+ configurations with diverse noise settings. Beyond final returns, we introduce the denoising factor to quantify the encoder's ability to filter distractions. To further isolate the effect of metric learning, we propose an isolated metric estimation setting, where the encoder is influenced solely by the metric loss. Our results show that metric learning improves return and denoising only marginally, as its benefits fade when key design choices, such as layer normalization and self-prediction loss, are incorporated into the baseline. We also find that commonly used benchmarks (e.g., grayscale videos, varying state-based Gaussian noise dimensions) add little difficulty, while Gaussian noise with random projection and pixel-based Gaussian noise remain challenging even for the best methods. Finally, we release an open-source, modular codebase to improve reproducibility and support future research on metric learning in deep RL.

# **Contribution(s)**

 We analyze five metric learning approaches under the isometric embedding framework to identify key design choices.
 Context: Metric learning methods often diverge significantly between theory and imple-

mentation.

- We introduce the denoising factor to quantify an encoder's ability to filter distractions.
   Context: Metric learning is often motivated by denoising ability but is rarely evaluated directly, with prior work relying mainly on qualitative analysis (Zhang et al., 2020).
- We benchmark five metric learning approaches across diverse distracting domains and find that common benchmarks add little difficulty to clean tasks, while certain noise settings remain challenging even for the best methods.
   Context: Prior work primarily uses IID Gaussian noise with varied dimensions (Ni et al., 2024) and grayscale video backgrounds (Zhang et al., 2020).
- Through ablation studies, we identify layer normalization and self-prediction loss as key design choices across all methods.
   Context: Prior work in metric learning does not isolate the effect of self-prediction loss and only shows the benefits of normalization in specific methods (Zang et al., 2022).
- We show that the benefits of metric learning diminish in both return and denoising factor when key design choices are incorporated into the baseline.
   Context: Prior work does not report this limitation of metric learning.

# **Understanding Effectiveness of Learning Behavioral Metrics in Deep Reinforcement Learning**

Anonymous authors

Paper under double-blind review

# Abstract

A key approach to state abstraction is approximating behavioral metrics (notably, bisim-1 2 ulation metrics) in the observation space, and embed these learned distances in the rep-3 resentation space. While promising for robustness to task-irrelevant noise shown in prior work, accurately estimating these metrics remains challenging, requiring various 4 5 design choices that create gaps between theory and practice. Prior evaluations focus 6 mainly on final returns, leaving the quality of learned metrics and the source of perfor-7 mance gains unclear. To systematically assess how metric learning works in deep RL, 8 we evaluate five recent approaches. We unify them under isometric embedding, identify 9 key design choices, and benchmark them with baselines across 20 state-based and 14 10 pixel-based tasks, spanning 250+ configurations with diverse noise settings. Beyond 11 final returns, we introduce the denoising factor to quantify the encoder's ability to fil-12 ter distractions. To further isolate the effect of metric learning, we propose an isolated 13 metric estimation setting, where the encoder is influenced solely by the metric loss.

14 Our results show that metric learning improves return and denoising only marginally, 15 as its benefits fade when key design choices, such as layer normalization and self-16 prediction loss, are incorporated into the baseline. We also find that commonly used benchmarks (e.g., grayscale videos, varying state-based Gaussian noise dimensions) 17 18 add little difficulty, while Gaussian noise with random projection and pixel-based Gaus-19 sian noise remain challenging even for the best methods. Finally, we release an opensource, modular codebase to improve reproducibility and support future research on 20 21 metric learning in deep RL.<sup>1</sup>

## 22 1 Introduction

Real-world environments often present high-dimensional, noisy observations, posing challenges for 23 24 RL. For instance, in image-based settings, task-irrelevant variations in background, lighting, and viewpoint introduce distractions. Yet, despite this observational complexity, system dynamics are 25 typically governed by a compact, task-relevant state. State abstraction (Li et al., 2006; Konidaris, 26 27 2019) provides a framework for extracting such *latent representations* from raw observations, fil-28 tering out irrelevant information while preserving task-critical structure. A key principle of state 29 abstraction is that behaviorally similar states should have similar representations. Traditionally, this 30 is enforced through state aggregation (Singh et al., 1994; Givan et al., 2003), grouping states into 31 discrete abstract classes based on equivalence relations. However, state aggregation lacks a measure 32 of *how different* states are across classes and struggles with continuous representations, requiring 33 infinitely many discrete classes.

34 To address this, bisimulation metrics (Ferns et al., 2004; 2011) and their scalable variants (Castro,

35 2020; Zhang et al., 2020) have been proposed to define meaningful distances between observations.

36 These fall under the broader class of **behavioral metrics** (Castro et al., 2023), which quantify state

<sup>&</sup>lt;sup>1</sup>The artifact is available at https://anonymous.4open.science/r/understanding\_metric-3C44

37 similarity based on differences in immediate rewards and transition probabilities. By learning a met-

ric alongside deep RL, prior work (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021;

39 Chen & Pan, 2022; Zang et al., 2022) has shown progress in tackling high-dimensional, noisy tasks.

40 Nevertheless, the role of behavioral metric learning in deep RL (metric learning for short) remains

41 unclear due to the lack of systematic evaluation. First, its effectiveness relies on accurately estimat-

42 ing these metrics, which is challenging in complex tasks. However, prior work primarily measures

43 performance through *returns*, without directly assessing the quality of the learned metrics. Second, 44 metric learning is often combined with *multiple losses* (e.g., self-prediction (Zhang et al., 2020), in-

44 incure rearing is often combined with *maniple tosses* (e.g., sen-prediction (Zhang et al., 2020), in-45 verse dynamics (Kemertas & Aumentado-Armstrong, 2021)), as well as *architectural choices* (e.g.,

46 normalization, ensembles (Zang et al., 2022)), making it difficult to isolate the contribution of metric

47 learning to performance gains. Third, most studies evaluate only OOD generalization in environ-

48 ments with grayscale natural videos as distractions (Zhang et al., 2020), conflating robustness with

49 generalization. Lastly, prior evaluations (Tomar et al., 2021; Li et al., 2022) report inconsistent

50 results for the same algorithms, raising concerns about reproducibility.

In this paper, we provide a understanding of how metric learning works in deep RL through 51 52 a systematic evaluation of five recent approaches alongside two baselines. First, we unify these 53 metric learning objectives under a common framework using the notion of isometric embedding, 54 identifying key design choices for our investigation. Next, to ensure a rigorous and comprehensive 55 evaluation, we introduce diverse distraction benchmarks by varying difficulty levels, from Gaussian 56 noise to colored natural videos, across both state-based and pixel-based domains, tested under both 57 ID and OOD generalization. Then, we quantify the **denoising** capability – the encoder's ability to 58 filter out distractions. We introduce the *denoising factor*, which numerically measures how well the 59 encoder distinguishes similar from dissimilar observations.<sup>2</sup> Finally, we propose an *isolated* metric 60 estimation setting to assess metric learning's contribution to denoising, independent of other losses.

61 **Contributions.** Our main contributions are as follows:

Conceptual: We analyze five recent metric learning approaches using the isometric embedding
 framework to identify key design choices. In addition, we propose *denoising factor* to quantify
 an agent's denoising capability in distracting tasks.

 Comprehensive benchmarking: We benchmark these five metric learning methods and baselines on diverse distracting variants of the DeepMind Control (DMC) suite (Tassa et al., 2018).
 In state-based domains, across 20 DMC tasks with 10 IID Gaussian noise settings, SimSR (Zang et al., 2022), originally evaluated in pixel-based domains, significantly outperforms other methods in both return and denoising factor. In pixel-based domains, across 14 DMC tasks with 6 image background settings, RAP (Chen & Pan, 2022) performs generally best. SAC (Haarnoja et al.,

2018) and DeepMDP (Gelada et al., 2019) remain competitive baselines but are often overlooked.

Reevaluating benchmark difficulty: Surprisingly, common distracting benchmarks – varying
 Gaussian noise dimensions in state-based domains and grayscale videos in pixel-based domains
 - add little difficulty. However, Gaussian noise with random projection in state-based domains
 and Gaussian noise in pixel-based domains remain challenging, even for the best methods.

4. Identifying key design choices: We find self-prediction loss is crucial to SimSR's success.
 Notably, (layer) normalization, used in SimSR, consistently improves return and denoising across all metric learning methods and baselines.

5. Marginal impact of metric learning (a bitter lesson): The benefits of metric learning diminish
 when key design choices are incorporated into the baseline.

81 6. **Open-source codebase**: We open-source a modular and efficient codebase to improve 82 reproducibility in the RL community.

<sup>&</sup>lt;sup>2</sup>Prior work (Zhang et al., 2020) qualitatively analyzes denoising by visualizing representations with t-SNE.

# 83 2 Background

## 84 2.1 Problem Formulation

85 We consider a setting where observations contain distractions and focus on a special class of Markov

86 decision processes – exogenous block MDPs (EX-BMDPs) (Efroni et al., 2021; Islam et al., 2022).

87 Before introducing EX-BMDPs, we first define block MDPs as a prerequisite.

88 **Block MDPs (Du et al., 2019).** A block MDP (BMDP) is a tuple  $\langle \mathcal{X}, \mathcal{Z}, \mathcal{A}, q, p, R, \gamma \rangle$ , where  $\mathcal{X}$ is the observation space, Z is the latent state space, A is the action space,  $p: Z \times A \to \Delta(Z)$  is 89 latent transition function,  $R: \mathbb{Z} \times \mathcal{A} \to \mathbb{R}$  is (latent) reward function, and  $\gamma \in [0, 1)$  is the discount 90 factor. The emission function  $q: \mathcal{Z} \to \Delta(\mathcal{X})$  generates observation  $x \sim q(\cdot \mid z)$  from latent 91 92 state z. Crucially, BMDP assumes the block structure:  $\forall z_1, z_2 \in \mathbb{Z}, z_1 \neq z_2 \implies \operatorname{supp}(q(\cdot \mid z_1, z_2 \in \mathbb{Z}))$  $(z_1) \cap \operatorname{supp}(q(\cdot \mid z_2)) = \emptyset$ . This ensures that each observation uniquely determines its latent state, 93 enabling the existence of the oracle encoder  $q^{-1}$ :  $\mathcal{X} \to \mathcal{Z}$  such that  $q^{-1}(x) = z$  whenever 94  $x \sim q(\cdot \mid z)$ . The goal of RL in BMDP is to find a policy  $\pi : \mathcal{X} \to \Delta(\mathcal{A})$  that maximizes the 95 rewards:  $\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} R(q^{-1}(x_{t}), a_{t}) \right]$ . The policy only receives the observation x without 96 access to the latent state z, the latent space  $\mathcal{Z}$ , or the oracle encoder  $q^{-1}$ . While the class of BMDPs 97 is equivalent to the class of MDPs (Du et al., 2019), they capture the underlying state from a high-98 99 dimensional observation. However, BMDPs do not differentiate between task-relevant (endogenous) 100 state and task-irrelevant (exogenous) noise in the latent space.

Exogenous BMDPs (Efroni et al., 2021). An EX-BMDP extends BMDP by factorizing a latent 101 state into  $z = (s, \xi)$ , where  $s \in S$  is the *task-relevant state* and  $\xi \in \Xi$  is the *task-irrelevant noise*, 102 representing distraction. The latent state transition  $p(s', \xi' \mid s, \xi, a)$  factorizes as  $p(s' \mid s, a)p(\xi' \mid \xi)$ , 103 where the noise  $\xi$  evolves independently and does not affect the reward function. To simplify nota-104 105 tion, we denote the reward function as as R(s, a). EX-BMDPs guarantee the existence of a denoising map  $D: \mathcal{Z} \to \mathcal{S}$  extracts the task-relevant state s from latent state  $z \in \mathcal{Z}$ . Combined with the oracle 106 107 encoder in BMDPs, this enables recovery of the task-relevant state directly from observations: s = $D(q^{-1}(x))$ . We define this composite function  $\phi^* = D \circ q^{-1}$  as the **oracle encoder** of EX-BMDP. 108

## 109 2.2 Representation Learning in RL

In actor-critic methods (Konda & Tsitsiklis, 1999), representation learning is commonly used to 110 111 handle complex MDPs such as EX-BMDPs. The idea is to learn an encoder that maps a raw obser-112 vation to a representation, which is then shared by both actor and critic. Formally, an actor-critic algorithm employs an encoder  $\phi: \mathcal{X} \to \Psi$ , a (latent) actor  $\pi_{\theta}: \Psi \to \Delta(\mathcal{A})$ , and a (latent) critic 113 114  $Q_{\omega}: \Psi \times \mathcal{A} \to \mathbb{R}$ , where  $\Psi$  is the representation space. In this work, we focus on end-to-end actorcritic methods based on the soft actor-critic (SAC) algorithm (Haarnoja et al., 2018). These methods 115 116 jointly optimize the encoder and actor-critic using the RL loss in SAC, denoted as  $J_{SAC}(\phi, \theta, \omega)$ . 117 Learning state representations solely from reward signals (i.e., RL loss) is challenging in complex 118 tasks. To address this, various state abstraction frameworks and representation objectives have 119 been proposed (see Ni et al. (2024) for a literature review). Among these, model-irrelevance abstraction (Li et al., 2006) defines two conditions for an effective encoder using bisimulation 120 121 relation (Givan et al., 2003). The first condition, known as **reward prediction** (**RP**)<sup>3</sup>, requires that

the representation preserves reward information. The second condition, known as **self-prediction**  $(\mathbb{ZP})^4$  (Ni et al., 2024), requires that the representation preserves latent dynamics information.

124 Model-irrelevance abstraction thus defines compact yet informative encoders that retain sufficient

125 information for optimal decision-making (Subramanian et al., 2022). By definition, the RP and ZP

<sup>126</sup> conditions hold when  $\phi = \phi^*$  and  $\Psi = S^5$ . This implies that the oracle encoder  $\phi^*$  serves as a 127 model-irrelevance abstraction.

<sup>&</sup>lt;sup>3</sup>Formally, in an EX-BMDP, RP condition is  $\exists R_{\kappa} : \Psi \times \mathcal{A} \to \mathbb{R}$ , s.t.  $R(\phi^*(x), a) = R_{\kappa}(\phi(x), a), \forall x, a$ .

<sup>&</sup>lt;sup>4</sup>Formally, in an EX-BMDP, ZP condition is  $\exists P_{\nu} : \Psi \times \mathcal{A} \to \Delta(\Psi)$ , s.t.  $P(\psi' \mid x, a) = P_{\nu}(\psi' \mid \phi(x), a), \forall x, a, \psi'$ . <sup>5</sup>In this case,  $R_{\kappa}(s, a) = R(s, a)$  and  $P_{\nu}(s' \mid s, a) = p(s' \mid s, a)$ .

- 128 To learn a model-irrelevance abstraction, **DeepMDP** (Gelada et al., 2019) introduces **RP and ZP**
- losses to approximate the RP and ZP conditions, respectively. Given a data tuple (x, a, r, x'), these
- 130 losses jointly optimizes the encoder  $\phi$ , the reward model  $R_{\kappa}$ , and the latent transition model  $P_{\nu}$ :

$$J_{\rm RP}(\phi,\kappa) = (R_{\kappa}(\phi(x),a) - r)^2, \quad J_{\rm ZP}(\phi,\nu) = -\log P_{\nu}(\bar{\phi}(x') \mid \phi(x),a), \tag{1}$$

where  $\overline{\phi}$  detaches the encoder from gradient back-propagation. The overall objective  $J_{\text{DeepMDP}}(\phi)$ for the encoder in DeepMDP combines SAC loss with RP and ZP losses (Eq. 1).

# 133 **3** Conceptual Analysis on Behavioral Metrics Learning in RL

This section establishes a conceptual framework linking behavioral metrics to representations in deep RL (Sec. 3.1), and then summarizes how related work instantiates it (Sec. 3.2).

## 136 3.1 Isometric Embedding: Between Behavioral Metrics and Representation

We aim to find an encoder that maps noisy observations into a structured representation space, where distances reflect differences in rewards and transition dynamics smoothly. This representation should facilitate RL by ensuring that task-relevant variations are captured. A natural way to formalize this goal is through the concept of an *isometric ambadding (isometry)*<sup>6</sup>:

formalize this goal is through the concept of an *isometric embedding* (*isometry*)<sup>6</sup>:

141 **Definition 1 (Isometric Embedding)** An encoder  $\phi : \mathcal{X} \to \Psi$  is an isometric embedding if the 142 distances in the original space  $(\mathcal{X}, d_{\mathcal{X}})$  are preserved in the representation space  $(\Psi, d_{\Psi})$ . Formally,

$$d_{\mathcal{X}}(x_1, x_2) = d_{\Psi}(\phi(x_1), \phi(x_2)), \quad \forall x_1, x_2 \in \mathcal{X},$$
(2)

143 where  $d_{\mathcal{X}}$  is the **target metric** ("desired" metric) and  $d_{\Psi}$  is the **representational metric**. See Ap-

144 *pendix Sec. B.1 for background on metric definitions.* 

## 145 3.2 Design Choices in Metric Learning

Table 1: Summary of key implementation choices for the benchmarked methods.

Method	$d_R$	$d_{\Psi}$	$d_T$	Other Losses	Transition Model	Metric Loss	Normalization	Target Trick
SAC	_	_	_	_	_	_	_	_
DeepMDP	_	_	_	RP + ZP	Probabilistic	_	_	_
DBC	Huber	Huber	$W_2$ closed-form	RP + ZP	Probabilistic	MSE	_	_
RDBC	Huber	Huber	$W_2$ closed-form	RP + ZP	Deterministic	MSE	Max norm	_
MICo	Abs.	MICo Angular	Sample-based	_	_	Huber		$\checkmark$
SimSR	Abs.	Cosine	Sample-based	ZP	Probabilistic Ensemble	Huber	$L_2$ norm	_
RAP	RAP	MICo Angular	$W_2$ closed-form	RP + ZP	Probabilistic	Huber	_	_

146 Def. 1 provides a general conceptual framework instantiated by several works in deep RL through 147 distinct design choices. Rather than delving into theoretical implications, we focus on *practical* 

*implementations* reflected in their publicly available codebases, summarized in Table 1.

149 **Choices of Target Metric**  $d_{\mathcal{X}}$ . The target metric, which captures differences in rewards and tran-150 sition dynamics, is typically formulated as (Castro et al., 2023): for  $x_1, x_2 \in \mathcal{X}$ ,

$$d_{\mathcal{X}}(x_1, x_2) = c_R d_R(r_1, r_2) + c_T d_T(d_{\mathcal{X}}) (P(\cdot \mid x_1), P(\cdot \mid x_2)), \tag{3}$$

- 151 which is inspired by bisimulation metrics (Ferns et al., 2004; 2011).<sup>7</sup>
- 152  $d_R$ , representing *immediate* state similarity, measures the closeness of sampled immediate rewards
- 153  $r_1, r_2 \in \mathbb{R}$  based on  $x_1, x_2$ . Common choices include absolute difference  $d_R(r_1, r_2) = |r_1 r_2|$
- 154 ("Abs." in Table 1) (Zang et al., 2022; Castro et al., 2021), "Huber" distance<sup>8</sup>  $d_R(r_1, r_2) =$

155 
$$\frac{1}{2}(r_1 - r_2)^2 \mathbf{1}_{\{|r_1 - r_2| < 1\}} + \left(|r_1 - r_2| - \frac{1}{2}\right) \mathbf{1}_{\{|r_1 - r_2| \ge 1\}}$$
 (Huber in Table 1), or other specific  
156 forms (Chen & Pan, 2022).

<sup>&</sup>lt;sup>6</sup>https://en.wikipedia.org/wiki/Isometry

<sup>&</sup>lt;sup>7</sup>See Appendix Sec. B.2 for the full formal definition and discussion of the variants of bisimulation metrics. <sup>8</sup>https://pytorch.org/docs/stable/generated/torch.nn.SmoothLlLoss.html

*d<sub>T</sub>*, a probabilistic measure of *long-term* state similarity, typically avoids expensive 1-Wasserstein
 computations in bisimulation metrics. Methods such as (Zhang et al., 2020; Kemertas &
 Aumentado-Armstrong, 2021; Chen & Pan, 2022) approximate transitions using a Gaussian tran sition model with a 2-Wasserstein metric. In contrast, Castro et al. (2021); Zang et al. (2022) rely
 on sample-based distance approximations.

162 **Choices of Representational Metric**  $d_{\Psi}$ . To approximate  $d_{\Psi}$ , methods employ either a Huber dis-163 tance (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021) (a surrogate for  $L_p$ -distance) 164 or an angular distance (Castro et al., 2021; Zang et al., 2022; Chen & Pan, 2022).

165 Metric Loss Function  $J_M$ . To approximate an isometric embedding, metric learning methods 166 optimize this general objective:

$$J_M(\phi) = \ell(d_{\Psi}(\phi(x_1), \phi(x_2)) - d_{\mathcal{X}}(x_1, x_2)), \tag{4}$$

where  $\ell$  can be Huber loss (Chen & Pan, 2022; Castro et al., 2021; Zang et al., 2022), or mean square error (MSE) (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021).

169 Self-prediction (ZP). As discussed, approximating  $d_{\mathcal{X}}$  often requires a transition model, methods 170 adopt distinct approaches: probabilistic models (Zhang et al., 2020; Castro et al., 2021), ensembles 171 of probabilistic models (Zang et al., 2022), and deterministic models (Kemertas & Aumentado-172 Armstrong, 2021). MICo, in contrast, employs a sample-based target metric that is free of ZP.

173 **Normalization.** We discuss normalization in the representation space  $\Psi$ . RDBC (Kemertas & 174 Aumentado-Armstrong, 2021) employs max normalization to enforce boundedness, leveraging prior 175 knowledge of target metric constraints. SimSR (Zang et al., 2022) applies  $L_2$  normalization to 176 enforce unit-length representations. LayerNorm (Ba et al., 2016) is widely used in pixel-based 177 encoder across all methods except SimSR. In the broader context of deep RL, normalization has 178 been extensively studied for its benefits in stabilizing training and generalization (Li et al., 2023; 179 Fujimoto et al., 2023; Elsayed et al., 2024; Gallici et al., 2024).

180 **Target trick.** MICo employs a target network  $\overline{\phi}$  for encoding one observation in  $d_{\Psi}$  when approx-181 imating  $d_{\chi}$  to ensure learning stability. See Castro et al. (2021, Appendix C.2) for further details.

## 182 3.3 Candidate Methods

183 We present the design choices of methods to be benchmarked in our work in Table 1. Note that 184 RDBC (Kemertas & Aumentado-Armstrong, 2021) incorporates additional components – such as 185 intrinsic rewards and inverse dynamics – that enhance performance. However, since these elements 186 are orthogonal to our study, they are not included in our implementation. For a fair comparison, our 187 experiments employ a probabilistic transition model for all methods that require one.

# 188 4 Study Design: Does Metric Learning Help with Denoising?

The "denoising capability" of behavioral metric learning is often cited as a motivation in prior work (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021; Chen & Pan, 2022; Zang et al., 2022). However, most studies evaluate this *indirectly* by (1) combining metric learning with RL, (2) training only on grayscale natural video backgrounds, (3) testing on unseen videos in training, and (4) evaluating solely through return performance. This leaves a gap between motivation and actual denoising assessment.

This section bridges that gap with a systematic study design. First, we introduce a diverse range of noise settings from IID Gaussian noise and random projections to natural video backgrounds (Sec. 4.1), enabling an analysis of how noise difficulty impacts metric learning. Second, we separate the noise distributions during training and testing to examine denoising under both ID and OOD



Figure 1: Examples of different noise settings in pixel-based domains (three consecutive timesteps each). From left to right: original clean image, IID Gaussian image, grayscale natural image, colored natural image, grayscale natural video, colored natural video. Three background images are different for video settings.

199 generalization settings (Sec. 4.2). Third, we introduce a direct evaluation measure, the *denoising* 

200 factor (Sec. 4.3). Finally, to disentangle metric learning from RL, we propose the isolated metric

201 estimation setting, where metric learning affects only the encoder, not the RL agent (Sec. 4.4).

## 202 4.1 Noise Settings

We describe our noise settings using the EX-BMDP framework (Sec. 2.1), where observations follow  $x \sim q(\cdot | z)$  with  $z = (s, \xi)$ .

**IID Gaussian Noise.** The task-irrelevant noise  $\xi_t$  is sampled independently at each timestep from an *m*-dimensional isotropic Gaussian,  $\xi_t \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ . For state-based domains, the observation is exactly the latent state, i.e.,  $x_t = z_t$  with *q* as the identity mapping. We adjust the noise dimension *m* or noise std  $\sigma$  to modulate difficulty, whereas prior work (Kemertas & Aumentado-Armstrong, 2021; Ni et al., 2024) only varies *m* with a small  $\sigma$ . For pixel-based domains, noise is applied per pixel in the background and overlaid by the robot foreground's pixels, with *q* as a rendering function.

111 **IID Gaussian Noise with Random Projection.** This setting applies only to state-based domains where  $s \in \mathbb{R}^n$ . Before interaction with the MDP, we construct a full-rank square matrix  $\mathbf{A} \in \mathbb{R}^{(n+m)\times(n+m)}$  with entries sampled as  $A_{ij} \sim \mathcal{N}(\mu_A, \sigma_A^2)$ . At each time step, we generate m-dimensional IID Gaussian noise  $\xi_t \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$  and then apply a linear projection to obtain observation  $x_t = \mathbf{A}z_t$  where  $z_t = (s_t, \xi_t)$ . Since  $\mathbf{A}$  is full rank,  $s_t$  can be recovered from  $x_t$  using  $\mathbf{A}^{-1}$ . This setting is more challenging than IID Gaussian noise, as it linearly entangles  $s_t$  and  $\xi_t$ , with q as the linear projection.<sup>9</sup>

**Natural Images.** This setting applies only to pixel-based domains, replacing the clean background with a randomly selected natural image. As in the original environment, the background remains fixed during training. Images can be *grayscale* or *colored*, introducing different levels of visual complexity. In EX-BMDP notation,  $\xi_t$  is stationary and q is a rendering function.

Natural Videos. This setting also applies only to pixel-based domains, replacing the clean background with with natural video. The underlying noise  $\xi_t \in \mathbb{N}$ , representing the *frame index*, follows the update rule  $\xi_t = (\xi_{t-1} + 1) \mod N$ , where N is the total number of frames. These videos can be *grayscale* or *colored*, with the grayscale setting widely used in metric learning (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021; Zang et al., 2022; Chen & Pan, 2022).

## 227 4.2 Denoising Involves ID and OOD Generalization

In this work, "*denoising*" refers to a form of generalization that removes task-irrelevant noise from observations, enabling generalization to tasks with unseen noise. The evaluation settings differ based on whether the *noise distribution* remains unchanged or shifts between training and testing.

<sup>&</sup>lt;sup>9</sup>Voelcker et al. (2024) similarly projects observations using a random binary matrix in discrete domains.

231 In-distribution (ID) Generalization Evaluation. The training and testing environments (EX-

BMDPs) are identical, meaning the same noise distribution is applied in both phases. For example,

233 IID Gaussian noise remains unchanged throughout training and testing.

Out-of-distribution (OOD) Generalization Evaluation. The training and testing EX-BMDPs share the same task-relevant parts (i.e., p(s' | s, a),  $p(s_0)$ , R(s, a)) but differ in noise distributions (i.e.,  $p(\xi' | \xi)$ ,  $p(\xi_0)$ ). For instance, natural videos from a training dataset are employed during training, while videos from a distinct test dataset are used during evaluation. This evaluation setup is widely used in metric learning (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021; Zang et al., 2022; Chen & Pan, 2022).

#### 240 4.3 Quantifying Denoising via the Denoising Factor

We introduce the *denoising factor* (DF), a measure that quantifies an encoder  $\phi$ 's ability to filter out irrelevant details while retaining essential information.<sup>10</sup> It also provides insight into how the behavioral metrics are approximated, given that exact behavioral metrics are nearly inaccessible via fixed-point iteration in high-dimensional state or action spaces. To compute DF, we define a *positive score* and a *negative score* for an encoder  $\phi$ . Inspired by triplet loss (Schroff et al., 2015) in contrastive learning, we compute these scores by selecting an observation x as an *anchor*, then constructing a *positive example*  $x^+$  similar to x, and a *negative example*  $x^-$  dissimilar to x.

**Definition 2 (Positive score)** The positive score of an encoder  $\phi$  w.r.t. the metric  $d_{\Phi}$  measures the average representational distance between anchors and their positive examples:

$$\operatorname{Pos}_{d_{\Phi}}(\phi) := \mathbb{E}_{x \sim \rho_{\pi}(x), \xi^{+} \sim \rho(\xi^{+}), x^{+} \sim q(\cdot | \phi^{*}(x), \xi^{+})} \left[ d_{\Phi}(\phi(x), \phi(x^{+})) \right], \tag{5}$$

where  $\rho_{\pi}(x)$  is the stationary state distribution under the policy  $\pi$  and  $\rho(\xi^+)$  is a stationary noise distribution. The sampling  $x^+ \sim q(\cdot \mid \phi^*(x), \xi^+)$  ensures that  $x^+$  shares the same task-relevant state  $s = \phi^*(x)$  but has different noise  $\xi^+$ .

In the temporally-independent noise setting,  $\rho(\xi^+)$  matches the noise transition; in the natural-video setting,  $\rho(\xi^+)$  is a uniform distribution over frame indices  $\{0, 1, \dots, N-1\}$ .

**Definition 3 (Negative score)** The negative score of an encoder  $\phi$  w.r.t. the metric  $d_{\Phi}$  measures the average representational distance between anchors and their negative examples:

$$\operatorname{Neg}_{d_{\Phi}}(\phi) := \mathbb{E}_{x, x^{-\operatorname{iid}} \sim \rho_{\pi}} \left[ d_{\Phi}(\phi(x), \phi(x^{-})) \right], \tag{6}$$

(7)

257 where  $x, x^-$  are IID samples from  $\rho_{\pi}$ .

**Definition 4 (Denoising factor (DF))** The denoising factor of an encoder  $\phi$  w.r.t. the metric  $d_{\Phi}$  is defined as the normalized difference between the negative and positive scores. As a result, a higher DF indicates better denoising ability.

$$\mathrm{DF}_{d_{\Phi}}(\phi) := \frac{\mathrm{Neg}_{d_{\Phi}}(\phi) - \mathrm{Pos}_{d_{\Phi}}(\phi)}{\mathrm{Neg}_{d_{\Phi}}(\phi) + \mathrm{Pos}_{d_{\Phi}}(\phi)} \in [-1, 1].$$

## 261 4.4 Decoupling Behavioral Metric Learning from RL

262 In many behavioral metric learning methods, the encoder  $\phi$  is optimized via a combination of losses: the RL loss (e.g.,  $J_{SAC}(\phi)$ ), the reward-prediction loss  $J_{RP}(\phi)$ , the self-prediction loss  $J_{ZP}(\phi)$ 263 264 (Eq. 1), and a metric loss  $J_{\rm M}(\phi)$  (3). This coupling makes it difficult to isolate the direct impact of 265 metric learning on representation quality. Moreover, denoising factor (DF, Def. 4) depends on both 266 the encoder and the data distribution induced by RL agent. Policies that frequently revisit similar 267 task-relevant states under varying noise may inflate DF, making it an unreliable measure of denois-268 ing. Due to the above reasons, we propose to evaluate behavioral metric learning algorithms in an 269 isolated metric estimation setting.

<sup>&</sup>lt;sup>10</sup>While the oracle encoder  $\phi^*$  achieves perfect denoising, direct comparison is impossible as  $\phi$  lacks access to S.

270 **Isolated Metric Estimation Setting.** To isolate the effect of metric learning, we introduce an iso-271 lated *metric encoder*  $\phi$  that is optimized solely via the metric loss  $J_{\rm M}(\phi)$ , while the *agent encoder*  $\phi$  is updated using the remaining training objectives (e.g.,  $J_{SAC}(\phi)$  or  $J_{DeepMDP}(\phi)$ ). In our exper-272 273 iments, regardless of the metric learning method, a SAC agent interacts with the environment and collect data for learning the metrics. This allows for a fair comparison of  $DF_{d,\phi}(\phi)$  across different 274 275 metric learning methods. For methods that rely on self-prediction loss (Zhang et al., 2020; Kemertas & Aumentado-Armstrong, 2021; Zang et al., 2022), we learn an *isolated transition model* using  $\phi$ 276 277 while preventing gradient backpropagation to  $\phi$  to ensure isolation.

# 278 5 Experiments

279 **Experiment Organization.** We first conduct a comprehensive evaluation of all the methods (Ta-280 ble 1) across 20 state-based DeepMind Control (DMC) (Tassa et al., 2018; Tunyasuvunakool et al., 281 2020) tasks (listed in Table 5) and 14 pixel-based DMC tasks (listed in Table 6). Evaluations are 282 performed under various noise settings using ID generalization. This study covers a significantly 283 larger task set than prior works. Our results (Sec. 5.1) provide a broad assessment of agent per-284 formance and task difficulty, as reflected by return variations. Based on these findings, we select 285 a subset of representative tasks for fine-grained analysis (Sec. 5.2) to examine the impact of key 286 design choices (Sec. 3.2). We further investigate the isolated metric evaluation setting (Sec. 4.4) in Sec. 5.3, and assess OOD generalization following prior work in Sec. 5.4. 287

**Evaluation Protocol.** We report the *mean episodic reward* rather than the IQM (Agarwal et al., 2021) to avoid ignoring tasks that are too easy or too challenging. For each run, the reported mean episodic reward, bounded within [0, 1000] for all tasks, is the average of 10 evaluation points within a 1.95M-2.05M step window and aggregated over seeds. Figures and tables display 95% confidence intervals over tasks. For state-based environments, we use 12 random seeds per *configuration*, where a configuration is defined as a (*task, noise setting*) pair. For pixel-based experiments, we use 5 random seeds per configuration.

Approximation of Denoising Factor (Eq. 7). All observations collected in the evaluation stage are regarded as *anchors*. We sample 16 positive samples and negative samples for each anchor using the strategy shown in Sec. 4.3. For consistency, we report  $DF_{\|\cdot\|_2}(\phi)$ .

## 298 5.1 Benchmarking Methods on Various Noise Settings

299 **Settings.** For state-based DMC tasks, we apply IID Gaussian noise, varying either (a) standard de-300 viations  $\sigma \in \{0.2, 1.0, 2.0, 4.0, 8.0\}$  (with a fixed noise dimension m = 32), or (b) noise dimensions 301  $m \in \{2, 16, 32, 64, 128\}$  (with a fixed standard deviation  $\sigma = 1.0$ ). For pixel-based DMC tasks, 302 evaluation is conducted under 6 image background settings: (1) clean background (the original 303 pixel-based DMC setting), (2) grayscale natural images, (3) colored natural images, (4) grayscale 304 natural videos, (5) colored natural videos, and (6) IID Gaussian noise (with  $\sigma = 1.0$ ). ID generaliza-305 tion evaluation is conducted in this subsection. The aggregated reward and DF for settings (a), (b), 306 and (1)-(6) are shown in Fig. 2 and Fig. 3, respectively. Per-task results are listed in Appendix Sec. E.

307 **Implementation Details.** For state-based tasks, the encoder is a three-layered MLP, as used by 308 SAC (Haarnoja et al., 2018) and RDBC (Kemertas & Aumentado-Armstrong, 2021). For pixel-309 based tasks, the encoder is a CNN followed by LayerNorm (Ba et al., 2016), as used by SAC-310 AE (Yarats et al., 2021b). All the compared methods are implemented based on SAC. For a fair 311 comparison, we adopt an identical probabilistic latent transition models and reward models used in 312 DBC and RDBC, although some methods, such as SimSR (Zang et al., 2022), propose using an 313 ensemble of latent transition models. We follow the exact hyperparameters specific to each metric 314 learning method. Please see Appendix Sec. D for further details.

315 Benchmarking Findings. We summarize the key findings from Fig. 2 and Fig. 3.



Figure 2: Benchmarking results: performance of seven methods across diverse noise settings, aggregating episodic rewards from 20 state-based (first two rows) and 14 pixel-based tasks (last row). "Noise std" denotes the IID Gaussian noise's standard deviation  $\sigma$ , while "noise dim" denotes its dimension m. Bars show 95% CI.



Figure 3: **Benchmarking results: reward (left) and denoising factor (right)** of seven methods to IID Gaussian noise dimension (Noise Dim) and standard deviation (Noise Std). Each point is aggregated by 20 state-based tasks in Table 5.

- SimSR consistently achieves the highest performance in most state-based tasks, excelling in both
   return and denoising factor. RAP performs best in most pixel-based tasks but suffers a moderate
   performance drop in state-based tasks. Interestingly, both SimSR and RAP were evaluated only
   in pixel-based domains in their papers, making our state-based findings novel.
- **SAC and DeepMDP**, as fundamental metric learning baselines, deliver decent performance on both pixel-based and state-based tasks. However, they are often overlooked in prior work.
- In state-based domains, all methods are generally more robust to increased noise *dimensions* (when  $\sigma = 1.0$ ), as commonly used in prior work, than to increased *standard deviation* (when m = 32), as shown in Fig. 3.
- In pixel-based domains, *grayscale natural video*, widely used in prior work, is not significantly
   harder than clean background setting (e.g., for SAC and DeepMDP). Surprisingly, the *IID Gaussian noise* setting is the most challenging, warranting further study.
- Different algorithms excel in different tasks (Table 5, Table 6), e.g., RAP in reacher/easy, MICo in point\_mass/easy (Fig. 23). Broad task coverage is essential to ensure generalizable insights.
- Adding objectives trades-off computation efficiency. As shown in Table 2, the time cost of optimizing of a metric loss is close to optimizing a ZP loss by comparing MICo with DeepMDP.

#### 332 5.2 What Matters in Metric and Representation Learning?

To identify key factors in metric learning in state-based domains, we conduct case studies on the design choices outlined in Sec. 3.2. We select *three medium-to-hard DMC tasks*, finger/turn\_easy, figure/turn\_hard, quadruped/run for detailed analysis. First, a notable difference between our default encoder implementations for state-based and pixel-based tasks is the inclusion of normalization, which may significantly impact benchmarking outcomes. SimSR, the best-performing algorithm in state-based environments, employs  $L_2$  normalization in the representation space and discusses its

Table 2: Relative time spent on model updates on NVIDIA L40S GPUs under the same task (walker/walk, with  $S = \mathbb{R}^{24}$  and  $\Xi = \mathbb{R}^{32}$ ). Values represent the multiple of SAC's updating time. Key hyperparameters affecting the speed are set identically to Table 7.



Figure 4: Case study on three DMC state-based tasks (IID Gaussian noise, noise dim=32, noise std=8.0), examining the effects of **including LayerNorm** (left vs. right three columns), **applying the target trick** (RDBC (T)), and **using Huber loss** (RDBC (H)) instead of MSE as the metric loss.

- effectiveness (Zang et al., 2022). This inspires us to examine whether normalization benefits *other* metric learning methods. Second, several techniques used by the best-performing methods merit
- 341 further analysis. For instance, SimSR, RAP, and MICo (which excels in colored natural video set-

342 tings) utilize Huber metric loss instead of MSE, while MICo incorporates the target trick (Sec. 3.2).

- 343 Third, we investigate the performance of methods with LayerNorm in a challenging setting: IID
- Gaussian noise with random projection (Sec. 4.1) with  $\sigma \in \{0.2, 2.0, 4.0, 8.0\}$  (with a fixed noise
- dimension m = 32), shown in Fig. 6. Important findings from Fig. 4 to Fig. 6 is as follows:
- Most methods benefit from LayerNorm in the representation space, improving both reward and DF (Fig. 4). Notably, DeepMDP with LayerNorm performs comparably to SimSR.<sup>11</sup> For RDBC, using Huber loss for the metric and using the target trick enhance performance (Fig. 4).<sup>12</sup>
- **ZP loss is crucial for SimSR's success** in noisy state-based tasks (Fig. 5).

A significant performance drop occurs for all agents when increasing the noise standard deviation
 in **IID Gaussian with random projection** setting (Fig. 6), even with normalization applied.
 Nevertheless, DeepMDP and SimSR remain relatively robust to the noise.

## 353 5.3 Isolated Metric Evaluation Setting: Does Learned Metrics Denoise?

- We further analyze the proposed setting in Sec. 4.4 in both state-based (shown in Fig. 7) and pixelbased settings (shown from Fig. 25 to Fig. 29). We observe that:
- Generally, metrics learned in isolation achieve some denoising but underperform compared to
   those co-learned with ZP and critic losses or even with ZP and critic losses alone in DeepMDP
   (first two rows of Fig. 7).
- LayerNorm also works well with isolated metric estimation (last two rows of Fig. 7).
- MICo's DF remains relatively low, which aligns with its theoretical property that the metric for positive samples is non-zero (Fig. 7), as MICo does not enforce zero self-distance.

<sup>&</sup>lt;sup>11</sup>Our additional experiments reveal that removing LayerNorm in pixel-based environments causes a substantial performance drop across all methods, highlighting the critical role of normalization.

<sup>&</sup>lt;sup>12</sup>We observe that runs with superior design choices exhibit much more stable representation norms and gradient norms for both the critic loss and metric loss. Thus, the performance gains shown in Fig. 4 are likely due to improved numerical stability in extrapolating the metrics and Q values during generalization.



Figure 5: Ablation study on ZP loss on SimSR. "SimSR (Original)" is the configuration benchmarked in Sec. 5.1, where ZP is integral to the metric estimation. Therefore, we resort to "SimSR (Basic)" setting (Theorem 2, Zang et al. (2022)), where the ZP component is independent from the metric estimation, and "SimSR (Basic, No ZP)" is the setting that ZP is detached from SimSR (Basic). This ablation highlights the impact of detaching ZP on the overall performance.



Figure 6: Aggregated reward (top row) and DF (bottom row) of seven agents on various **IID Gaussian with** random projection settings in the 3 selected state-based tasks.

#### 362 **5.4 OOD Generalization**

While prior work has focused on OOD generalization in pixel-based settings, we extend this analysis by evaluating all **14 pixel-based tasks**. The "grayscale video" setting (and similarly for other settings) in Fig. 8 denotes using grayscale videos for both training and evaluation, with distinct video samples in each phase. Takeaways in Fig. 8 and Fig. 9 are as follows:

Comparing Fig. 8 (OOD) with Fig. 2 (ID), all methods struggle to generalize in both grayscale and colored image settings. Unlike video backgrounds, which provide temporal variation, static images lack diverse cues, making adaptation to unseen backgrounds more challenging.

- Even with OOD generalization evaluation, SAC and DeepMDP remain competitive baselines (Fig. 8).
- OOD generalization is more challenging in the *colored* video setting than in the *grayscale* video setting (Fig. 9). Surprisingly, even baselines like SAC exhibit a low reward gap, questioning the
- necessity of incorporating a metric loss in the widely-used grayscale setting.



Figure 8: Aggregated reward (top row) and DF (bottom row) of seven agents on various noise settings in 14 pixel-based tasks in Table 6 with **OOD generalization** evaluation.

Under review for RLC 2025, to be published in RLJ 2025



Figure 7: Top row: DF for the agent encoder  $\phi$  (co-trained with RL in Sec. 5.1) without LayerNorm. Mid row: DF for the isolated metric encoder  $\tilde{\phi}$  without LayerNorm. Bottom row: DF for  $\tilde{\phi}$  with LayerNorm. All use ID generalization evaluation.



Figure 9: **Reward gap** (performance in ID evaluation minus OOD evaluation) in the grayscale video setting (left) and the colored video setting (right), aggregated on 14 pixel-based tasks in Table 6.

# 375 6 Conclusion

Takeaways. Based on our empirical studies on behavioral metric learning in deep RL, we highlight the following key insights for RL researchers:

- To gain a clearer understanding of RL algorithms, initial evaluations should be conducted on simple, controlled environments (e.g., varying Gaussian noise std, pixel-based Gaussian noise).
- Claims and motivations for metric learning should be supported by rigorous evaluation, including
   measures such as the denoising factor and comparisons between ID and OOD generalization.
- 382 3. Self-prediction loss and LayerNorm are critical design choices that significantly impact metricand representation learning.
- The benefits of metric learning diminish when key design choices, such as self-prediction loss
   and LayerNorm, are integrated into SAC. This calls for a deeper investigation into when and how
   metric learning provides unique advantages beyond these existing techniques.

Future work. Our study focuses on continuous control; future work should explore discrete domains and real-world tasks. The relationship between denoising and return performance remains unclear, requiring further investigation. Additionally, improved benchmarks are needed to better isolate the effects of metric learning.

# 391 References

- 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 informa-*
- *tion processing systems*, 34:29304–29320, 2021.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Pablo Samuel Castro. Scalable methods for computing state similarity in deterministic markov
   decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34,
   pp. 10069–10076, 2020.
- Pablo Samuel Castro, Tyler Kastner, Prakash Panangaden, and Mark Rowland. Mico: Improved
   representations via sampling-based state similarity for markov decision processes. *Advances in Neural Information Processing Systems*, 34:30113–30126, 2021.
- Pablo Samuel Castro, Tyler Kastner, Prakash Panangaden, and Mark Rowland. A kernel perspective
   on behavioural metrics for markov decision processes. *arXiv preprint arXiv:2310.19804*, 2023.
- Jianda Chen and Sinno Pan. Learning representations via a robust behavioral metric for deep re inforcement learning. Advances in Neural Information Processing Systems, 35:36654–36666,
   2022.
- Edmund M Clarke. Model checking. In *Foundations of Software Technology and Theoretical Com- puter Science: 17th Conference Kharagpur, India, December 18–20, 1997 Proceedings 17*, pp.
  54–56. Springer, 1997.
- Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.
  Provably efficient rl with rich observations via latent state decoding. In *International Conference on Machine Learning*, pp. 1665–1674. PMLR, 2019.
- Yonathan Efroni, Dipendra Misra, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Provable rl with exogenous distractors via multistep inverse dynamics. *arXiv preprint arXiv:2110.08847*, 2021.
- Mohamed Elsayed, Gautham Vasan, and A Rupam Mahmood. Streaming deep reinforcement learning finally works. *arXiv preprint arXiv:2410.14606*, 2024.
- Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite markov decision processes.
   In *Proceedings of the 20th conference on Uncertainty in artificial intelligence*, pp. 162–169, 2004.
- Norm Ferns, Prakash Panangaden, and Doina Precup. Bisimulation metrics for continuous markov
   decision processes. *SIAM Journal on Computing*, 40(6):1662–1714, 2011.
- Scott Fujimoto, Wei-Di Chang, Edward Smith, Shixiang Shane Gu, Doina Precup, and David Meger.
   For sale: State-action representation learning for deep reinforcement learning. *Advances in neural information processing systems*, 36:61573–61624, 2023.
- Matteo Gallici, Mattie Fellows, Benjamin Ellis, Bartomeu Pou, Ivan Masmitja, Jakob Nicolaus
  Foerster, and Mario Martin. Simplifying deep temporal difference learning. *arXiv preprint arXiv:2407.04811*, 2024.
- Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp:
   Learning continuous latent space models for representation learning. In *International conference on machine learning*, pp. 2170–2179. PMLR, 2019.
- Robert Givan, Thomas Dean, and Matthew Greig. Equivalence notions and model minimization in
  markov decision processes. *Artificial intelligence*, 147(1-2):163–223, 2003.

- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
  maximum entropy deep reinforcement learning with a stochastic actor. In *International confer- ence on machine learning*, pp. 1861–1870. Pmlr, 2018.
- Riashat Islam, Manan Tomar, Alex Lamb, Yonathan Efroni, Hongyu Zang, Aniket Didolkar, Dipendra Misra, Xin Li, Harm Van Seijen, Remi Tachet des Combes, et al. Agent-controller represen-
- tations: Principled offline rl with rich exogenous information. *arXiv preprint arXiv:2211.00164*,
  2022.
- Mete Kemertas and Tristan Aumentado-Armstrong. Towards robust bisimulation metric learning.
   *Advances in Neural Information Processing Systems*, 34:4764–4777, 2021.
- Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- George Konidaris. On the necessity of abstraction. *Current opinion in behavioral sciences*, 29:1–7,
  2019.
- Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstraction
  for mdps. *AI&M*, 1(2):3, 2006.
- Lu Li, Jiafei Lyu, Guozheng Ma, Zilin Wang, Zhenjie Yang, Xiu Li, and Zhiheng Li. Normalization
  enhances generalization in visual reinforcement learning. *arXiv preprint arXiv:2306.00656*, 2023.
- 451 Xiang Li, Jinghuan Shang, Srijan Das, and Michael Ryoo. Does self-supervised learning really im452 prove reinforcement learning from pixels? *Advances in Neural Information Processing Systems*,
  453 35:30865–30881, 2022.
- Tianwei Ni, Benjamin Eysenbach, Erfan Seyedsalehi, Michel Ma, Clement Gehring, Aditya Mahajan, and Pierre-Luc Bacon. Bridging state and history representations: Understanding selfpredictive rl. arXiv preprint arXiv:2401.08898, 2024.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face
   recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- 460 Satinder Singh, Tommi Jaakkola, and Michael Jordan. Reinforcement learning with soft state ag-461 gregation. *Advances in neural information processing systems*, 7, 1994.
- Jayakumar Subramanian, Amit Sinha, Raihan Seraj, and Aditya Mahajan. Approximate information
   state for approximate planning and reinforcement learning in partially observed systems. *Journal of Machine Learning Research*, 23(12):1–83, 2022.
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- Manan Tomar, Utkarsh A Mishra, Amy Zhang, and Matthew E Taylor. Learning representations for
   pixel-based control: What matters and why? *arXiv preprint arXiv:2111.07775*, 2021.
- 470 Saran Tunyasuvunakool, Alistair Muldal, Yotam Doron, Siqi Liu, Steven Bohez, Josh Merel, Tom
  471 Erez, Timothy Lillicrap, Nicolas Heess, and Yuval Tassa. dm\_control: Software and tasks for
  472 continuous control. *Software Impacts*, 6:100022, 2020.
- 473 Claas Voelcker, Tyler Kastner, Igor Gilitschenski, and Amir-massoud Farahmand. When does
  474 self-prediction help? understanding auxiliary tasks in reinforcement learning. *arXiv preprint*475 *arXiv:2406.17718*, 2024.
- Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous con trol: Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021a.

- 478 Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improv-
- ing sample efficiency in model-free reinforcement learning from images. In *Proceedings of the*
- 480 *aaai conference on artificial intelligence*, volume 35, pp. 10674–10681, 2021b.
- 481 Hongyu Zang, Xin Li, and Mingzhong Wang. Simsr: Simple distance-based state representations
- for deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*,
  volume 36, pp. 8997–9005, 2022.
- 484 Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning
- invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020.

487	Supplementary Materials
488 489	The following content was not necessarily subject to peer review.

## 490 A Notation

491 Table 3 and Table 4 show the glossary used in this paper.

Table 3: Glossary of notations in EX-BMDP (Sec. 2.1). The top section lists symbols related to the latent states, while the bottom section defines symbols related to grounded observations.

Symbol	Description
$z = (s,\xi) \in \mathcal{Z}$	Environment's latent state
$p(z' \mid z, a)$	Latent state transition
$s \in \mathcal{S}$	Task-relevant state
$\xi \in \Xi$	Task-irrelevant noise
$a \in \mathcal{A}$	Action
R(s, a)	Latent reward function
$r \in \mathbb{R}$	Reward
$\gamma \in [0,1)$	Discount factor
$p(s' \mid s, a)$	Task-relevant state transition
$p(\xi' \mid \xi)$	Task-irrelevant noise transition
$x \in \mathcal{X}$	Observation
$q(x \mid z)$	Emission function
$q^{-1}: \mathcal{X} \to \mathcal{Z}$	Oracle encoder to $\mathcal{Z}$
$\phi^*:\mathcal{X}\to\mathcal{S}$	Oracle encoder (to $S$ )
$\mathcal{R}(x,a)$	Grounded reward function
$\mathcal{P}(x' \mid x, a)$	Grounded transition function

Table 4: Glossary of notations in RL agents.

Symbol	Description
$\psi\in \Psi$	Agent's representation
$\phi: \mathcal{X} \to \Psi$	Agent's encoder
$\pi_{\theta}: \Psi \to \Delta(\mathcal{A})$	(Latent) Actor
$Q_{\omega}: \Psi \times \mathcal{A} \to \mathbb{R}$	(Latent) Critic
$R_{\kappa}:\Psi\times\mathcal{A}\to\mathbb{R}$	(Latent) Reward model
$P_{\nu}: \Psi \times \mathcal{A} \to \Delta(\Psi)$	(Latent) Transition model
$\tilde{\phi}: \mathcal{X} \to \Psi$	(Isolated) Metric encoder

#### 492 **B** Background on Metrics

#### 493 B.1 Metric, Pseudometric, and Diffuse Metric

494 **Metric.** A function  $d : \mathcal{X} \times \mathcal{X} \to \mathbb{R}_{\geq 0}$  is called a *metric* on the space  $\mathcal{X}$  if for all  $x, y, z \in \mathcal{X}$ :

(1) 
$$d(x,y) = 0 \iff x = y,$$
  
(2)  $d(x,y) = d(y,x),$   
(3)  $d(x,z) \le d(x,y) + d(y,z).$ 

495 **Pseudometric**<sup>13</sup>. A function  $d : \mathcal{X} \times \mathcal{X} \to \mathbb{R}_{\geq 0}$  is a *pseudometric* if it satisfies (2) and (3) above, 496 and for (1) we only require d(x, x) = 0 for all  $x \in \mathcal{X}$  (i.e., d(x, y) = 0 does not imply x = y).

497 **Diffuse Metric (Definition 4.9 in Castro et al. (2021)).** A function  $d : \mathcal{X} \times \mathcal{X} \to \mathbb{R}_{\geq 0}$  is a *diffuse* 498 *metric* if it satisfies properties (2) and (3) above, and for (1) we only require  $d(x, y) \ge 0$ . That is, 499 we do not demand d(x, x) = 0 or that  $d(x, y) = 0 \iff x = y$ .

#### 500 B.2 Definitions of Various Behavioral Metrics

501 In this section, we discuss the behavioral metrics for an EX-BMDP (Sec. 2.1). From the observation

502 space  $\mathcal{X}$ , the grounded transition function is defined as  $\mathcal{P}(x' \mid x, a) = \sum_{z' \in \mathcal{Z}} q(x' \mid z') p(z' \mid z')$ 

<sup>&</sup>lt;sup>13</sup>Sometimes termed "semimetric" (Ferns et al., 2004).

503  $q^{-1}(x), a$  and the grounded reward function as  $\mathcal{R}(x, a) = R(\phi^*(x), a)$ . Let  $x_i, x_j \in \mathcal{X}$  be two 504 arbitrary observations.

Bisimulation metric is a relaxation of bisimulation relation (Givan et al., 2003) by allowing a smooth
variation based on differences in the reward function and transition dynamics. Originally rooted in
the concept of identical future reward sequences, it is closely linked to model checking (Clarke,
1997). The bisimulation metric thus quantifies the behavioral similarity between two states and is
formally defined as follows:

510 **Definition 5 (Bisimulation metric**  $d^{\sim}$  (Ferns et al., 2004; 2011)) There exists a unique pseudo-511 metric  $d^{\sim} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ , called the bisimulation metric (BSM)<sup>14</sup>, defined as:

$$d^{\sim}(x_i, x_j) := \max_{a \in \mathcal{A}} \left( c_R |\mathcal{R}(x_i, a) - \mathcal{R}(x_j, a)| + c_T \mathcal{W}_1(d^{\sim}) (\mathcal{P}(\cdot \mid x_i, a), \mathcal{P}(\cdot \mid x_j, a)) \right), \quad (8)$$

512 where *I*-Wasserstein distance  $W_1(d^{\sim})(P,Q) = \inf_{\delta} \mathbb{E}_{(x'_i,x'_j)\sim\delta} [d^{\sim}(x'_i,x'_j)]$ , the coupling  $\delta$  satifies 513  $\sum_{x'_j} \delta(x'_i,x'_j) = P(x'_i)$ , and  $\sum_{x'_i} \delta(x'_i,x'_j) = Q(x'_j)$ ,  $c_R$  and  $c_T$  are coefficients for short-term and 514 long-term behavioral differences.

515 In practice, applying the max operator over actions is intractable in continuous action spaces and

516 pessimistically accounts for behavioral similarity across all actions, including those leading to low 517 rewards. Policy-dependent bisimulation metrics (Castro, 2020) address this limitation by restricting 518 behavioral similarity to the current policy, eliminating the need to evaluate all actions.

519 **Definition 6 (Policy-dependent bisimulation metric**  $d^{\pi}$  (Castro, 2020)) *Given a policy*  $\pi : \mathcal{X} \to \Delta(\mathcal{A})$ , there exists a unique pseudometric  $d^{\pi} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ , called a policy-dependent bisimulation *metric (PBSM), defined as:* 

$$d^{\pi}(x_i, x_j) := c_R |\mathcal{R}^{\pi}(x_i) - \mathcal{R}^{\pi}(x_j)| + c_T \mathcal{W}_1(d^{\pi})(\mathcal{P}^{\pi}(\cdot \mid x_i), \mathcal{P}^{\pi}(\cdot \mid x_j)),$$

(9)

522 where  $\mathcal{R}^{\pi}(x) := \mathbb{E}_{a \sim \pi}[\mathcal{R}(x, a)]$  and  $\mathcal{P}^{\pi}(\cdot \mid x) := \mathbb{E}_{a \sim \pi}[\mathcal{P}^{\pi}(\cdot \mid x, a)]$  are policy-dependent reward 523 and transition, respectively.

To approximate the 1-Wasserstein distance in PBSM, DBC (Zhang et al., 2020) assumes a Gaussian transition kernel and uses 2-Wasserstein distance which has a close-form solution under such assumption. To further circumvent the costly computation of the 1-Wasserstein distance, Castro et al. (2021) proposes MICo distance.

528 **Definition 7 (MICo distance**  $u^{\pi}$  (**Castro et al., 2021**)) Given a policy  $\pi : \mathcal{X} \to \Delta(\mathcal{A})$ , there exists 529 a unique diffuse metric  $u^{\pi} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ , called MICo distance:

$$u^{\pi}(x_{i}, x_{j}) := c_{R} |\mathcal{R}^{\pi}(x_{i}) - \mathcal{R}^{\pi}(x_{j})| + c_{T} \mathbb{E}_{x_{i}^{\prime} \sim \mathcal{P}^{\pi}(\cdot|x_{i}), x_{j}^{\prime} \sim \mathcal{P}^{\pi}(\cdot|x_{j})} \left[ u^{\pi}(x_{i}^{\prime}, x_{j}^{\prime}) \right].$$
(10)

Based on MICo, several improvements are made by SimSR (Zang et al., 2022) and RAP (Chen &
Pan, 2022).

532 **Definition 8 (SimSR distance (Zang et al., 2022))** Given a policy  $\pi : \mathcal{X} \to \Delta(\mathcal{A})$ , there exists a 533 unique distance  $u^{\pi} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ , called the Simple State Representation (SimSR) distance:

$$u^{\pi}(x_{i}, x_{j}) := c_{R} \left| \mathcal{R}^{\pi}(x_{i}) - \mathcal{R}^{\pi}(x_{j}) \right| + c_{T} \mathbb{E}_{x_{i}^{\prime} \sim \widehat{\mathcal{P}}^{\pi}(\cdot|x_{i}), x_{j}^{\prime} \sim \widehat{\mathcal{P}}^{\pi}(\cdot|x_{j})} \left[ u^{\pi}(x_{i}^{\prime}, x_{j}^{\prime}) \right],$$
(11)

534 where  $u^{\pi}(x_i, x_j) = \cos(\phi(x_i), \phi(x_j))$  is the cosine distance and  $\widehat{\mathcal{P}}^{\pi}$  is an approximated transition 535 dynamics model.

536 **Definition 9 (RAP distance (Chen & Pan, 2022))** Given a policy  $\pi : \mathcal{X} \to \Delta(\mathcal{A})$ , there exists a 537 unique distance  $u^{\pi} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ , called the Robust Approximate (RAP) distance:

$$u^{\pi}(x_i, x_j) := c_R \left| \mathcal{R}^{\pi}(x_i) - \mathcal{R}^{\pi}(x_j) \right| + c_T \mathbb{E}_{a_i \sim \pi, a_j \sim \pi} \left[ u^{\pi} \left( \mathbb{E}_{x'_i \sim \hat{\mathcal{P}}(\cdot | x_i, a_i)}[x'_i], \mathbb{E}_{x'_j \sim \hat{\mathcal{P}}(\cdot | x_j, a_j)}[x'_j] \right) \right].$$

$$(12)$$

<sup>&</sup>lt;sup>14</sup>As noted by Ferns et al. (2004), BSM relates to the largest bisimulation relation,  $\sim$ . For brevity, we simplify the original definition that uses the fixed-point of an operator and omit the proof for existence of such a fixed-point.

# 538 C Difficulty Levels of the Tasks

Difficulty levels for each task in state-based DMC are provided in Table 5. These levels are based on the average reward across all compared methods. In comparison with the difficulty assignment in pixel-based DMC within DrQ-v2 (Yarats et al., 2021a, Table 1), most tasks labeled as easy in our evaluation are also easy in DrQ-v2, and most tasks labeled as medium or hard in our evaluation are medium in DrQ-v2. However, finger spin and pendulum swingup shift from easy to medium, and hopper stand from easy to hard.

Table 5: Difficulty levels for 20 state-based DMC tasks, as determined by the compared methods (Table 1). "Avg Reward" stands for the average reward across all IID Gaussian noise settings in Sec. 5.1. For each run, the reported reward is the average of 10 evaluation points collected around 2M steps. "Max (Min) Reward" denotes the best (worst) agent's average reward over 12 runs, while "Max/Min" is the ratio of the best to worst performance, indicating a task's ability to discriminate between agent performances.

Task		Avg Reward	Max Reward	Min Reward	Max/Min	Difficulty
ball_in_cup	catch	934.8	977.4	841.7	1.2	Easy
cartpole	balance	919.4	997.3	791.2	1.3	Easy
cartpole	balance_sparse	877.7	983.6	772.3	1.3	Easy
walker	stand	834.6	979.0	437.8	2.2	Easy
cartpole	swingup	818.1	874.1	707.6	1.2	Easy
walker	walk	805.7	961.9	382.4	2.5	Easy
reacher	easy	740.1	955.1	453.0	2.1	Medium
finger	spin	728.8	923.6	498.5	1.9	Medium
quadruped	walk	703.1	948.9	245.5	3.9	Medium
cartpole	swingup_sparse	647.3	839.1	531.9	1.6	Medium
reacher	hard	641.1	853.0	340.3	2.5	Medium
finger	turn_easy	587.8	926.5	207.7	4.5	Medium
walker	run	545.8	776.1	117.4	6.6	Medium
cheetah	run	533.4	859.0	129.8	6.6	Medium
pendulum	swingup	514.3	824.5	247.2	3.3	Medium
quadruped	run	460.7	864.3	199.0	4.3	Hard
finger	turn_hard	435.6	893.0	102.6	8.7	Hard
hopper	stand	261.9	878.4	22.3	39.3	Hard
acrobot	swingup	75.7	246.1	11.2	22.0	Hard
hopper	hop	64.7	243.4	1.5	162.4	Hard

Table 6: Difficulty levels for 14 pixel-based DMC tasks, as determined by the compared methods (Table 1). "Avg Reward" stands for the average reward across clean background, natural video (colored and grayscale), natural image (colored and grayscale), and IID Gaussian noise settings described in Sec. 5.1. For each run, the reported reward is the average of 10 evaluation points collected around 2M steps. "Max (Min) Reward" denotes the best (worst) agent's average reward over 5 runs, while "Max/Min" is the ratio of the best to worst performance, indicating a task's ability to discriminate between agent performances.

Task		Avg Reward	Max Reward	Min Reward	Max/Min	Difficulty
cartpole	balance	949.3	986.7	905.5	1.1	Easy
cartpole	balance_sparse	915.3	999.4	804.6	1.2	Easy
walker	stand	887.7	959.1	633.7	1.5	Easy
finger	spin	815.2	909.5	426.4	2.1	Easy
cartpole	swingup	765.0	853.2	551.1	1.5	Medium
ball_in_cup	catch	719.2	887.8	263.4	3.4	Medium
walker	walk	718.9	909.1	360.8	2.5	Medium
point_mass	easy	421.5	558.6	256.1	2.2	Medium
cartpole	swingup_sparse	409.6	680.4	57.8	11.8	Medium
reacher	easy	336.8	949.1	113.0	8.4	Hard
pendulum	swingup	313.2	468.9	9.3	50.5	Hard
cheetah	run	299.4	411.0	144.7	2.8	Hard
walker	run	285.5	441.9	77.3	5.7	Hard
hopper	hop	73.6	122.5	5.6	22.0	Hard

# 545 **D** Hyperparameters

546 We align our hyperparameter settings with the agents we benchmark, referring to their open-source 547 codebases and reported values in papers. Table 7 shows the general hyperparameter setting for most agents. For MICo (Castro et al., 2021)<sup>15</sup>,  $\beta$  in MICo distance parametrization is set to 0.1. For 548 RAP (Chen & Pan, 2022)<sup>16</sup>, we use the same hyperparameter setting in their open-source code, 549 where actor, critic, and encoder learning rate is set to  $5 \times 10^{-4}$ ,  $\beta$  in MICo distance parametrization 550 551 is set to  $10^{-6}$ , RP and ZP loss coefficients is set to  $10^{-4}$ , and encoder feature dimensionality is set 552 to 100. For action repeat, we set it to 4 for most tasks, to 8 for cartpole (swingup, swingup sparse), 553 and to 2 for finger spin and walker (walk, run, stand) following the convention (Yarats et al., 2021a; 554 Zang et al., 2022; Chen & Pan, 2022).

<sup>&</sup>lt;sup>15</sup>https://github.com/google-research/google-research/tree/master/mico <sup>16</sup>https://github.com/jianda-chen/RAP\_distance

Hyperparameter Name	Value
Replay buffer capacity	$1 \times 10^{6}$
Replay ratio	0.2
Batch size	128
Discount $\gamma$	0.99
Optimizer	Adam
Encoder feature dimensionality	50
Hidden units in neural networks	256
Paralleled environments	10
Critic learning rate	$1 \times 10^{-3}$
Critic target update frequency	2
Critic Q-function soft-update rate $\tau_Q$	0.01
Actor learning rate	$1 \times 10^{-3}$
Actor update frequency	2
Actor log stddev bounds	[-10, 2]
Encoder learning rate	$1 \times 10^{-3}$
Encoder soft-update rate $ au_{\phi}$	0.05
Reward model and transition model's learning rate	$1 \times 10^{-3}$
Reward model and transition model's weight decay	$1 \times 10^{-7}$
SAC temperature learning rate	$1 \times 10^{-4}$
SAC initial temperature	0.1
Metric loss weight $\lambda_M$	0.5
Metric reward coefficient $c_R$	1
Metric transition coefficient $c_T$	0.99
Image size	$84 \times 84 \times 3$
Frame stack	3
Distracting video frames $N$ (Per paralleled environment)	1000

Table 7: Hyperparameter settings for most agents we benchmarked.

# 555 E Additional Experiment Results

## 556 E.1 Per-task Result for State-based DMC with IID Gaussian Noise Settings

- 557 In the figures and tables shown in this section, the tasks are sorted by their difficulty levels shown in
- 558 Table 5 and 6.

Table 8: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=0.2, noise dim=32.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$958.0 \pm 3.9$	$971.0 \pm 4.0$	$883.5 \pm 114.0$	$975.0 \pm 1.8$	$974.2 \pm 1.4$	$966.6 \pm 8.2$	$\textbf{978.3} \pm 0.9$
cartpole	balance	$912.7 \pm 24.6$	$836.8 \pm 155.5$	$870.9 \pm 172.0$	$988.7 \pm 7.4$	$968.8 \pm 16.7$	$967.6 \pm 9.9$	$\textbf{995.0} \pm 5.9$
cartpole	balance_sparse	$790.2 \pm 143.2$	$579.4 \pm 282.6$	$760.3 \pm 193.4$	$624.0 \pm 217.8$	$815.9 \pm 182.7$	$882.7 \pm 167.7$	$\textbf{991.9} \pm 17.8$
walker	stand	$621.9 \pm 164.8$	$929.3 \pm 30.3$	$842.7 \pm 65.6$	$974.6 \pm 5.0$	$959.2 \pm 28.5$	$846.2\pm90.2$	$977.4 \pm 4.3$
cartpole	swingup	$776.0 \pm 95.9$	$790.1 \pm 143.1$	$872.7 \pm 2.3$	$582.5 \pm 203.8$	$867.0 \pm 3.9$	$839.4 \pm 13.0$	$\textbf{876.6} \pm 5.2$
walker	walk	$443.1\pm133.0$	$819.2\pm62.2$	$869.8 \pm 63.7$	$\textbf{963.8} \pm 3.9$	$947.2\pm 6.4$	$780.6 \pm 110.7$	$961.9 \pm 7.0$
reacher	easy	$474.1\pm219.0$	$501.3 \pm 262.0$	$768.0 \pm 169.8$	$556.0 \pm 221.6$	$861.5 \pm 146.6$	$831.6 \pm 153.3$	$\textbf{957.0} \pm 16.1$
finger	spin	$813.5 \pm 36.1$	$837.9 \pm 20.5$	$629.6 \pm 152.9$	$905.8 \pm 43.3$	$951.4 \pm 16.9$	$861.8 \pm 45.0$	$\textbf{980.8} \pm 4.0$
quadruped	walk	$209.4 \pm 55.9$	$820.5 \pm 53.7$	$890.3 \pm 47.4$	$778.1 \pm 135.0$	$614.2 \pm 234.7$	$817.7 \pm 98.6$	$\textbf{948.3} \pm 10.5$
cartpole	swingup_sparse	$615.1 \pm 188.8$	$729.4 \pm 148.4$	$800.5 \pm 32.9$	$661.6 \pm 196.7$	$824.0 \pm 14.1$	$575.9 \pm 221.7$	$\textbf{835.8} \pm 8.7$
reacher	hard	$675.8 \pm 115.8$	$530.3 \pm 237.7$	$322.6 \pm 188.7$	$436.5 \pm 234.7$	$\textbf{874.2} \pm 55.2$	$778.5 \pm 200.7$	$651.3 \pm 160.1$
finger	turn_easy	$231.1 \pm 44.4$	$759.3 \pm 81.1$	$660.5 \pm 140.5$	$806.7 \pm 58.3$	$853.2 \pm 124.2$	$772.3 \pm 120.8$	$\textbf{914.2} \pm 18.7$
walker	run	$219.5 \pm 76.2$	$464.8 \pm 28.6$	$730.4 \pm 12.1$	$747.0 \pm 8.0$	$668.0 \pm 12.9$	$368.3 \pm 123.1$	$\textbf{794.9} \pm 7.4$
cheetah	run	$134.3\pm107.7$	$462.9 \pm 61.3$	$529.4 \pm 102.4$	$668.4 \pm 11.9$	$673.8 \pm 23.0$	$667.2 \pm 47.0$	$\textbf{858.6} \pm 16.4$
pendulum	swingup	$389.6 \pm 220.7$	$774.3 \pm 128.6$	$165.3\pm204.8$	$836.4 \pm 5.4$	$380.8\pm255.2$	$364.2\pm265.2$	$\pmb{841.2} \pm 4.9$
quadruped	run	$176.5 \pm 41.0$	$446.3 \pm 20.4$	$542.9 \pm 39.3$	$462.3\pm55.0$	$360.3 \pm 71.1$	$435.6 \pm 21.6$	$\pmb{868.4} \pm 29.1$
finger	turn_hard	$102.2\pm20.6$	$588.7 \pm 158.8$	$242.6 \pm 152.8$	$580.2 \pm 133.3$	$774.9 \pm 192.2$	$846.1 \pm 35.6$	$\textbf{880.9} \pm 25.8$
hopper	stand	$11.0 \pm 9.4$	$255.0 \pm 150.3$	$36.8 \pm 14.0$	$411.6 \pm 160.4$	$219.7 \pm 61.4$	$517.2 \pm 193.2$	$\textbf{895.5} \pm 24.8$
acrobot	swingup	$33.2 \pm 11.6$	$\textbf{257.2} \pm 28.4$	$9.5 \pm 2.6$	$86.3 \pm 42.5$	$68.6 \pm 69.0$	$10.9 \pm 2.2$	$214.2\pm53.9$
hopper	hop	$0.8\pm0.2$	$5.3\pm 6.5$	$5.0\pm2.5$	$107.2\pm36.8$	$10.1 \pm 9.0$	$144.6\pm32.7$	$\textbf{252.4} \pm 8.0$

Table 9: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=1, noise dim=32.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$881.7 \pm 157.0$	$969.1 \pm 3.7$	$940.7 \pm 56.6$	$974.0 \pm 1.7$	$972.5 \pm 2.2$	$972.7 \pm 1.5$	<b>978.5</b> ± 0.8
cartpole	balance	$863.2 \pm 31.4$	$972.7 \pm 10.2$	$616.4 \pm 262.2$	$982.3 \pm 14.5$	$964.7 \pm 22.2$	$948.6 \pm 21.6$	$999.2 \pm 0.8$
cartpole	balance_sparse	$843.9 \pm 93.4$	$953.7 \pm 27.6$	$703.9 \pm 136.3$	$908.6 \pm 68.6$	$919.7 \pm 59.9$	$934.1 \pm 45.5$	$990.7 \pm 11.7$
walker	stand	$527.4 \pm 166.8$	$956.6 \pm 9.9$	$827.2 \pm 114.0$	$974.6 \pm 3.9$	$968.7 \pm 5.1$	$781.8 \pm 122.4$	$\textbf{977.1} \pm 3.6$
cartpole	swingup	$824.7 \pm 22.7$	$770.3 \pm 100.2$	$829.9 \pm 78.8$	$741.2 \pm 99.6$	$865.7 \pm 1.9$	$837.1 \pm 9.1$	$865.3 \pm 6.4$
walker	walk	$561.2 \pm 116.3$	$852.9\pm50.4$	$895.6 \pm 57.1$	$959.5 \pm 8.2$	$949.3 \pm 4.9$	$805.0 \pm 44.4$	$\textbf{959.6} \pm 13.7$
reacher	easy	$487.0 \pm 143.6$	$650.4 \pm 169.1$	$875.3 \pm 57.9$	$547.4 \pm 141.8$	$900.6 \pm 79.7$	$900.0 \pm 38.4$	$\textbf{959.5} \pm 9.3$
finger	spin	$626.2\pm202.6$	$826.6 \pm 23.4$	$519.5 \pm 170.5$	$904.8 \pm 32.4$	$897.9 \pm 22.1$	$851.5 \pm 52.3$	$981.4 \pm 3.6$
quadruped	walk	$314.9 \pm 151.6$	$753.4 \pm 53.4$	$875.7 \pm 61.5$	$922.6 \pm 51.8$	$534.5\pm201.9$	$752.1 \pm 157.3$	$948.0 \pm 10.0$
cartpole	swingup_sparse	$703.6 \pm 93.4$	$817.2 \pm 8.6$	$566.6 \pm 102.2$	$665.9 \pm 121.6$	$753.3 \pm 98.2$	$796.2 \pm 11.9$	$\textbf{836.6} \pm 9.4$
reacher	hard	$600.7 \pm 118.0$	$607.6 \pm 156.4$	$438.2 \pm 135.5$	$640.7 \pm 162.9$	$825.3 \pm 108.1$	$\textbf{882.5} \pm 71.6$	$864.5 \pm 131.0$
finger	turn_easy	$201.1 \pm 42.4$	$722.5 \pm 94.8$	$404.3 \pm 162.9$	$723.1 \pm 104.7$	$830.1 \pm 137.2$	$773.7 \pm 84.9$	$\textbf{928.4} \pm 17.6$
walker	run	$216.8 \pm 69.5$	$432.9 \pm 63.0$	$711.1 \pm 15.0$	$725.3 \pm 35.3$	$670.2 \pm 16.5$	$416.5 \pm 125.5$	$774.2 \pm 12.9$
cheetah	run	$21.3 \pm 45.0$	$509.8 \pm 21.5$	$435.9 \pm 129.0$	$653.0 \pm 22.0$	$646.1 \pm 15.5$	$707.5 \pm 23.8$	$863.2 \pm 21.1$
pendulum	swingup	$353.2\pm257.3$	$827.5 \pm 12.3$	$461.3 \pm 243.1$	$770.3 \pm 151.4$	$220.6\pm238.5$	$284.9 \pm 228.7$	$\textbf{833.8} \pm 12.6$
quadruped	run	$267.3 \pm 68.3$	$437.1 \pm 18.7$	$492.2 \pm 42.7$	$471.4 \pm 13.8$	$416.1\pm 81.2$	$447.8 \pm 16.5$	$\textbf{846.3} \pm 41.6$
finger	turn_hard	$105.1 \pm 17.2$	$603.1 \pm 112.9$	$174.1 \pm 61.9$	$728.5 \pm 70.8$	$588.6 \pm 178.1$	$513.3 \pm 131.8$	$908.7 \pm 13.9$
hopper	stand	$33.8 \pm 19.6$	$148.7\pm76.6$	$34.7 \pm 9.4$	$293.7 \pm 152.9$	$109.7 \pm 43.2$	$387.7 \pm 196.9$	$\textbf{848.0} \pm 112.5$
acrobot	swingup	$37.9 \pm 10.0$	$137.2\pm40.2$	$55.8 \pm 35.9$	$76.5 \pm 44.0$	$16.2\pm3.7$	$12.5 \pm 2.2$	$\textbf{263.5} \pm 51.9$
hopper	hop	$2.8 \pm 2.3$	$1.1 \pm 0.3$	$6.4 \pm 1.7$	$84.5 \pm 18.7$	$9.4 \pm 8.4$	$134.5 \pm 26.5$	$\textbf{245.8} \pm 13.3$

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$938.0 \pm 16.3$	$968.9 \pm 3.0$	$908.7 \pm 78.4$	$974.9 \pm 1.4$	$973.9 \pm 1.2$	$972.0 \pm 2.0$	$\textbf{978.2} \pm 0.6$
cartpole	balance	$892.5 \pm 44.6$	$956.9 \pm 28.5$	$768.2 \pm 242.0$	$982.4 \pm 11.8$	$962.1 \pm 21.2$	$948.8 \pm 25.3$	$\textbf{997.8} \pm 4.3$
cartpole	balance_sparse	$684.2 \pm 234.4$	$956.2 \pm 48.1$	$755.7 \pm 159.8$	$962.7 \pm 45.4$	$973.1 \pm 31.2$	$973.6 \pm 24.0$	$\textbf{990.5} \pm 17.5$
walker	stand	$416.2 \pm 166.5$	$930.5 \pm 32.5$	$853.4 \pm 89.4$	$969.5 \pm 15.1$	$965.2 \pm 9.5$	$763.8 \pm 116.7$	$\textbf{977.9} \pm 7.2$
cartpole	swingup	$711.8 \pm 144.2$	$852.4 \pm 7.9$	$864.4 \pm 9.6$	$523.6 \pm 199.3$	$854.6 \pm 6.6$	$829.4 \pm 11.2$	$\textbf{870.3} \pm 8.6$
walker	walk	$501.3 \pm 105.7$	$858.9 \pm 50.4$	$922.2\pm140.2$	$947.2 \pm 15.2$	$933.7 \pm 17.9$	$742.6 \pm 134.5$	$\textbf{958.3} \pm 10.3$
reacher	easy	$463.1\pm227.2$	$613.7 \pm 277.7$	$841.3 \pm 151.1$	$509.1 \pm 249.7$	$936.3 \pm 28.8$	$948.6 \pm 14.1$	$\textbf{957.7} \pm 22.8$
finger	spin	$521.1 \pm 214.8$	$811.3 \pm 15.6$	$512.7 \pm 134.7$	$902.1 \pm 26.2$	$850.5 \pm 13.1$	$873.2 \pm 36.1$	$\textbf{973.3} \pm 13.5$
quadruped	walk	$320.5 \pm 98.4$	$780.4 \pm 39.8$	$804.2\pm208.6$	$808.8 \pm 172.7$	$626.0 \pm 213.3$	$714.2 \pm 148.5$	$\textbf{954.2} \pm 3.9$
cartpole	swingup_sparse	$665.0 \pm 100.0$	$786.7 \pm 17.7$	$514.1 \pm 149.1$	$780.9 \pm 14.8$	$741.8 \pm 148.9$	$706.9 \pm 146.2$	$\pmb{840.9} \pm 6.4$
reacher	hard	$690.7 \pm 144.1$	$658.9 \pm 219.9$	$358.0 \pm 211.9$	$648.8 \pm 207.0$	$824.6 \pm 167.1$	$807.9 \pm 187.3$	$\textbf{923.8} \pm 32.9$
finger	turn_easy	$230.0 \pm 56.4$	$697.0 \pm 99.8$	$273.4 \pm 56.4$	$727.2 \pm 77.3$	$712.2 \pm 143.6$	$657.5 \pm 125.2$	$\textbf{935.0} \pm 14.5$
walker	run	$80.3 \pm 52.7$	$435.4 \pm 56.5$	$696.2 \pm 13.9$	$731.7 \pm 16.4$	$666.0 \pm 11.5$	$410.2 \pm 83.1$	$\textbf{777.6} \pm 10.4$
cheetah	run	$162.9 \pm 100.0$	$483.2 \pm 16.7$	$232.4 \pm 181.0$	$615.7 \pm 17.6$	$615.9 \pm 26.9$	$665.7 \pm 60.9$	$\textbf{851.5} \pm 15.2$
pendulum	swingup	$256.6 \pm 233.7$	$820.5 \pm 14.3$	$557.6 \pm 202.7$	$732.6 \pm 167.0$	$199.7\pm222.2$	$382.9 \pm 254.0$	$\textbf{835.1} \pm 9.6$
quadruped	run	$166.9 \pm 42.6$	$440.8 \pm 21.8$	$552.7 \pm 52.1$	$439.6 \pm 51.6$	$379.0 \pm 77.1$	$414.8 \pm 48.6$	$\pmb{864.5} \pm 40.2$
finger	turn_hard	$91.6 \pm 18.0$	$499.7 \pm 157.5$	$228.8 \pm 134.7$	$590.9 \pm 87.0$	$427.8 \pm 195.8$	$399.3 \pm 146.5$	$\pmb{897.2} \pm 19.5$
hopper	stand	$17.8 \pm 12.9$	$77.0 \pm 64.7$	$32.2 \pm 12.1$	$346.3 \pm 185.8$	$99.8 \pm 63.1$	$544.8 \pm 178.8$	$\textbf{904.7} \pm 20.4$
acrobot	swingup	$23.5 \pm 11.3$	$75.4 \pm 30.1$	$36.1 \pm 24.8$	$102.6 \pm 46.4$	$18.8 \pm 18.7$	$10.9\pm3.0$	$\textbf{244.1} \pm 47.0$
hopper	hop	$0.4 \pm 0.3$	$2.8 \pm 2.8$	$4.3 \pm 4.1$	$79.1 \pm 26.9$	$0.2\pm0.2$	$129.9 \pm 41.5$	$\textbf{242.6} \pm 15.6$

Table 10: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=2, noise dim=32.

Table 11: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=4, noise dim=32.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$918.6 \pm 41.8$	$967.0 \pm 3.5$	$723.9 \pm 139.3$	$972.2 \pm 4.3$	$973.9 \pm 1.8$	$960.3 \pm 15.7$	$\textbf{977.2} \pm 1.1$
cartpole	balance	$815.1 \pm 91.9$	$957.9 \pm 13.2$	$901.4 \pm 171.1$	$977.8 \pm 19.4$	$969.2 \pm 11.5$	$931.9 \pm 29.9$	$\textbf{996.4} \pm 5.7$
cartpole	balance_sparse	$831.1 \pm 158.1$	$917.9 \pm 50.4$	$830.1 \pm 131.0$	$932.6 \pm 103.1$	$996.2 \pm 3.6$	$899.8 \pm 86.1$	$991.1 \pm 17.2$
walker	stand	$394.7 \pm 180.1$	$929.3 \pm 28.9$	$875.1 \pm 57.7$	$971.1 \pm 6.4$	$951.9 \pm 19.4$	$641.9 \pm 146.9$	$\textbf{978.8} \pm 5.8$
cartpole	swingup	$770.4 \pm 34.2$	$850.4 \pm 5.4$	$856.4 \pm 21.9$	$684.2 \pm 176.0$	$858.5 \pm 5.6$	$823.9 \pm 15.7$	$\textbf{881.2} \pm 0.5$
walker	walk	$323.7 \pm 131.3$	$834.2\pm56.4$	$925.5\pm30.6$	$954.7 \pm 5.1$	$925.2 \pm 26.5$	$778.4 \pm 119.8$	$\textbf{965.8} \pm 4.5$
reacher	easy	$312.9 \pm 184.3$	$732.4 \pm 244.1$	$616.1 \pm 234.2$	$579.6 \pm 269.9$	$950.6 \pm 23.8$	$828.8 \pm 110.3$	$\textbf{951.3} \pm 18.1$
finger	spin	$220.8 \pm 190.4$	$719.2 \pm 68.4$	$528.9 \pm 134.6$	$917.6 \pm 22.1$	$681.0 \pm 144.3$	$862.5\pm50.0$	$875.1 \pm 46.6$
quadruped	walk	$204.0 \pm 65.6$	$700.7 \pm 105.9$	$768.2 \pm 206.0$	$709.7 \pm 147.2$	$506.2 \pm 209.1$	$878.9 \pm 27.7$	$\textbf{953.6} \pm 9.1$
cartpole	swingup_sparse	$468.4 \pm 219.8$	$562.2 \pm 180.6$	$459.9 \pm 198.6$	$617.6 \pm 182.8$	$65.2 \pm 143.5$	$385.4 \pm 256.0$	$\textbf{841.4} \pm 9.2$
reacher	hard	$369.3 \pm 239.1$	$615.5 \pm 221.1$	$389.7 \pm 203.8$	$624.0 \pm 230.1$	$580.1 \pm 280.0$	$923.0 \pm 32.1$	$942.7 \pm 12.6$
finger	turn_easy	$184.8 \pm 28.4$	$503.9 \pm 88.5$	$229.7 \pm 41.4$	$724.5\pm57.0$	$584.0 \pm 171.8$	$491.6 \pm 80.8$	$927.2 \pm 10.9$
walker	run	$57.6 \pm 46.2$	$462.2 \pm 43.1$	$674.3 \pm 15.2$	$730.9 \pm 19.1$	$646.4 \pm 20.1$	$432.1 \pm 118.4$	$766.5 \pm 13.0$
cheetah	run	$185.9 \pm 79.0$	$459.6 \pm 24.6$	$386.2 \pm 154.2$	$575.5 \pm 25.8$	$607.8 \pm 13.2$	$657.4 \pm 41.2$	$\textbf{857.6} \pm 20.8$
pendulum	swingup	$123.2\pm159.2$	$831.8 \pm 9.2$	$577.4 \pm 196.0$	$777.4 \pm 130.4$	$72.3 \pm 151.5$	$162.9 \pm 194.8$	$\textbf{838.7} \pm 4.0$
quadruped	run	$208.0 \pm 48.5$	$452.7 \pm 16.6$	$495.4 \pm 92.6$	$476.3 \pm 46.3$	$322.6 \pm 86.0$	$434.1 \pm 24.7$	$\textbf{850.7} \pm 25.8$
finger	turn_hard	$100.5\pm19.7$	$312.6 \pm 92.0$	$135.7\pm45.0$	$537.5 \pm 117.1$	$371.1 \pm 152.8$	$215.3 \pm 83.6$	$\textbf{904.8} \pm 20.6$
hopper	stand	$5.3\pm0.6$	$30.3 \pm 25.9$	$25.2 \pm 11.0$	$179.1\pm101.0$	$21.9 \pm 25.2$	$332.9 \pm 200.1$	$\textbf{828.0} \pm 114.3$
acrobot	swingup	$19.4 \pm 10.4$	$43.2 \pm 28.5$	$16.6 \pm 13.8$	$78.4 \pm 45.2$	$11.4 \pm 2.7$	$10.7\pm2.1$	$267.3 \pm 49.5$
hopper	hop	$0.2\pm0.2$	$0.7\pm0.3$	$1.3 \pm 1.0$	$75.0 \pm 22.7$	$0.1\pm0.0$	$51.8 \pm 28.0$	$\textbf{255.7} \pm 13.1$

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$924.7 \pm 20.2$	$965.2 \pm 3.9$	$760.4 \pm 142.3$	$969.3 \pm 4.8$	$924.8 \pm 95.5$	$957.7 \pm 22.7$	$\textbf{977.0} \pm 1.0$
cartpole	balance	$814.1 \pm 86.6$	$966.6 \pm 9.2$	$950.3 \pm 71.2$	$973.7 \pm 12.4$	$967.5 \pm 12.3$	$928.7 \pm 32.3$	$\textbf{999.5} \pm 0.5$
cartpole	balance_sparse	$646.6 \pm 200.8$	$935.6 \pm 34.8$	$784.0 \pm 139.7$	$982.3 \pm 29.3$	$984.7 \pm 11.4$	$823.5 \pm 160.2$	$\textbf{985.0} \pm 15.4$
walker	stand	$265.4 \pm 118.4$	$932.2 \pm 13.8$	$796.8 \pm 63.3$	$975.8 \pm 5.1$	$968.6 \pm 8.1$	$522.5 \pm 154.5$	$\textbf{978.5} \pm 5.5$
cartpole	swingup	$778.1 \pm 28.9$	$851.9 \pm 8.6$	$865.1 \pm 21.4$	$838.7 \pm 9.4$	$858.1 \pm 5.7$	$828.6 \pm 11.5$	$\textbf{878.8} \pm 4.9$
walker	walk	$115.8\pm60.1$	$834.5\pm78.2$	$817.5 \pm 111.6$	$954.6 \pm 8.8$	$919.9 \pm 23.6$	$537.1 \pm 133.7$	$\textbf{966.7} \pm 3.0$
reacher	easy	$306.0 \pm 135.1$	$585.7 \pm 274.6$	$552.9 \pm 232.2$	$605.0\pm201.1$	$896.6 \pm 70.9$	$906.6 \pm 22.3$	$\textbf{945.9} \pm 44.7$
finger	spin	$174.4 \pm 171.0$	$560.1\pm50.3$	$408.0 \pm 54.1$	$841.9 \pm 88.1$	$226.0 \pm 155.0$	$767.5 \pm 115.8$	$\textbf{849.3} \pm 53.8$
quadruped	walk	$255.5 \pm 75.2$	$777.8 \pm 61.2$	$805.9 \pm 134.0$	$804.2 \pm 113.3$	$407.6 \pm 237.5$	$690.5 \pm 142.0$	$\textbf{953.0} \pm 4.8$
cartpole	swingup_sparse	$526.6 \pm 169.5$	$526.0 \pm 202.3$	$199.2 \pm 150.1$	$613.2 \pm 181.5$	$0.0\pm0.0$	$318.3 \pm 250.6$	$\pmb{844.0} \pm 1.8$
reacher	hard	$395.4 \pm 245.3$	$367.7 \pm 231.4$	$157.7 \pm 100.0$	$578.9 \pm 163.3$	$630.9 \pm 247.0$	$808.4 \pm 173.4$	$\textbf{950.8} \pm 8.1$
finger	turn_easy	$201.9 \pm 38.5$	$419.0 \pm 75.9$	$240.6 \pm 36.4$	$619.0 \pm 35.1$	$592.9 \pm 176.6$	$327.3 \pm 88.5$	$\textbf{926.8} \pm 10.9$
walker	run	$23.9 \pm 2.6$	$455.9 \pm 41.3$	$649.4 \pm 11.1$	$628.9 \pm 25.7$	$635.3 \pm 19.8$	$347.8 \pm 84.0$	$\textbf{760.6} \pm 19.4$
cheetah	run	$185.6 \pm 86.7$	$433.8 \pm 34.2$	$335.0 \pm 87.5$	$553.2 \pm 26.4$	$578.8 \pm 22.0$	$628.9 \pm 49.6$	$\pmb{866.3} \pm 10.1$
pendulum	swingup	$124.1\pm155.1$	$753.1 \pm 148.5$	$462.9 \pm 222.3$	$736.8 \pm 243.9$	$4.0\pm3.2$	$10.9 \pm 2.7$	$\pmb{840.6} \pm 5.2$
quadruped	run	$219.5 \pm 63.5$	$417.9 \pm 44.2$	$441.1 \pm 93.7$	$433.3 \pm 47.3$	$233.8 \pm 59.0$	$381.1 \pm 64.9$	$\pmb{847.4} \pm 21.7$
finger	turn_hard	$97.9 \pm 11.8$	$207.2 \pm 53.8$	$110.8 \pm 17.0$	$414.7 \pm 49.5$	$177.6 \pm 66.1$	$168.3\pm50.4$	$\textbf{885.4} \pm 24.5$
hopper	stand	$5.8\pm0.5$	$29.1 \pm 21.1$	$25.1 \pm 14.4$	$198.1 \pm 103.0$	$11.4 \pm 10.7$	$116.6 \pm 45.1$	$\pmb{899.1} \pm 18.4$
acrobot	swingup	$11.1 \pm 5.8$	$19.8 \pm 10.9$	$17.0 \pm 9.8$	$72.5 \pm 43.9$	$14.3\pm5.6$	$9.8 \pm 3.6$	$\textbf{280.8} \pm 32.6$
hopper	hop	$0.3\pm0.3$	$0.4 \pm 0.3$	$0.8 \pm 0.5$	$51.1 \pm 13.4$	$0.1\pm0.0$	$31.3 \pm 16.7$	$\textbf{233.9} \pm 22.6$

Table 12: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=8, noise dim=32.

Table 13: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=1, noise dim=2.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$935.1 \pm 19.2$	$973.5 \pm 2.1$	$973.2 \pm 3.2$	$973.5 \pm 2.7$	$975.4 \pm 1.2$	$970.1 \pm 5.3$	$\textbf{978.0} \pm 1.0$
cartpole	balance	$894.1 \pm 36.9$	$648.2 \pm 248.1$	$934.6 \pm 56.1$	$983.9 \pm 22.2$	$950.5 \pm 18.7$	$935.2 \pm 27.7$	$\textbf{990.8} \pm 18.6$
cartpole	balance_sparse	$890.2 \pm 81.4$	$548.1 \pm 204.1$	$885.9 \pm 48.8$	$861.6 \pm 124.5$	$932.6 \pm 80.2$	$\textbf{976.8} \pm 19.6$	$949.6 \pm 85.3$
walker	stand	$737.8 \pm 123.8$	$951.9 \pm 9.4$	$907.9 \pm 31.3$	$970.3 \pm 8.3$	$957.6 \pm 17.1$	$903.4 \pm 38.6$	$976.6 \pm 9.3$
cartpole	swingup	$806.9 \pm 42.5$	$835.2 \pm 36.0$	$874.8 \pm 1.6$	$757.6 \pm 118.0$	$862.9 \pm 10.2$	$844.1 \pm 14.8$	$873.9 \pm 5.5$
walker	walk	$635.2 \pm 123.5$	$847.0\pm76.0$	$833.6 \pm 135.6$	$961.6 \pm 5.1$	$946.4 \pm 4.3$	$821.9 \pm 99.4$	$\textbf{964.7} \pm 4.2$
reacher	easy	$637.4 \pm 186.6$	$660.9 \pm 262.6$	$714.6 \pm 214.8$	$361.8 \pm 201.3$	$858.6 \pm 114.2$	$701.9 \pm 212.0$	$\textbf{929.2} \pm 18.6$
finger	spin	$789.3 \pm 41.8$	$875.8 \pm 28.6$	$440.7 \pm 151.8$	$922.9 \pm 33.5$	$969.0 \pm 6.6$	$903.1 \pm 29.2$	$979.5 \pm 7.1$
quadruped	walk	$270.3 \pm 131.4$	$787.2 \pm 56.5$	$851.2 \pm 96.3$	$783.7 \pm 150.7$	$760.2 \pm 178.3$	$807.6 \pm 129.3$	$950.5 \pm 11.0$
cartpole	swingup_sparse	$387.3 \pm 250.8$	$813.4 \pm 15.9$	$685.2 \pm 100.0$	$718.1 \pm 146.1$	$787.0 \pm 87.2$	$643.8 \pm 191.9$	$\textbf{834.8} \pm 9.8$
reacher	hard	$430.6 \pm 206.5$	$467.3 \pm 225.3$	$185.3\pm104.2$	$516.1 \pm 249.5$	$784.3 \pm 106.5$	$\textbf{837.9} \pm 127.6$	$592.1 \pm 165.2$
finger	turn_easy	$226.9 \pm 46.8$	$705.8 \pm 118.6$	$691.6 \pm 185.4$	$753.7 \pm 123.1$	$\textbf{934.9} \pm 10.0$	$829.5 \pm 98.3$	$912.3 \pm 16.3$
walker	run	$245.5 \pm 86.2$	$467.9 \pm 23.5$	$733.7 \pm 14.4$	$751.9 \pm 10.0$	$681.1 \pm 10.3$	$458.9 \pm 84.5$	$781.6 \pm 13.9$
cheetah	run	$167.6 \pm 98.0$	$493.7 \pm 45.9$	$484.8 \pm 106.5$	$662.6 \pm 24.8$	$699.0 \pm 13.2$	$688.6 \pm 52.0$	$\textbf{851.4} \pm 25.4$
pendulum	swingup	$513.3 \pm 232.5$	$\pmb{840.2} \pm 4.9$	$39.7\pm55.0$	$767.2 \pm 151.8$	$565.5\pm259.6$	$412.3\pm262.0$	$770.2\pm152.5$
quadruped	run	$170.7\pm52.5$	$448.3 \pm 13.0$	$571.1 \pm 60.6$	$472.0 \pm 59.4$	$356.0 \pm 90.1$	$396.8 \pm 59.3$	$\textbf{882.9} \pm 24.5$
finger	turn_hard	$131.0\pm72.2$	$702.8 \pm 118.4$	$329.1 \pm 223.8$	$516.0 \pm 175.0$	$689.5 \pm 198.0$	$755.5 \pm 113.1$	$\pmb{864.5} \pm 47.5$
hopper	stand	$22.8 \pm 18.7$	$199.8\pm140.6$	$39.0 \pm 13.0$	$363.7 \pm 180.1$	$218.6 \pm 59.1$	$351.6 \pm 145.5$	$\textbf{901.3} \pm 25.4$
acrobot	swingup	$25.6 \pm 7.7$	$\textbf{280.6} \pm 21.2$	$9.2\pm3.0$	$83.0 \pm 34.8$	$46.3 \pm 41.4$	$11.2 \pm 2.3$	$183.4 \pm 66.7$
hopper	hop	$1.5\pm2.0$	$2.4 \pm 2.2$	$4.4 \pm 2.3$	$83.2 \pm 29.2$	$52.6 \pm 39.9$	$110.0\pm45.2$	$\textbf{250.5} \pm 15.2$

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$956.9 \pm 4.9$	$971.9 \pm 2.3$	$966.5 \pm 10.6$	$975.5 \pm 1.3$	$974.4 \pm 1.8$	$966.5 \pm 14.1$	$978.1 \pm 0.8$
cartpole	balance	$865.9 \pm 56.5$	$932.1 \pm 58.2$	$768.3 \pm 210.6$	$990.9 \pm 5.0$	$958.0 \pm 31.8$	$933.1 \pm 17.9$	<b>996.4</b> ± 5.6
cartpole	balance_sparse	$764.4 \pm 120.6$	$820.0 \pm 134.4$	$830.0 \pm 79.4$	$727.8 \pm 135.0$	$962.2 \pm 31.3$	$937.8 \pm 38.4$	$987.1 \pm 12.8$
walker	stand	$514.7 \pm 122.9$	$940.2 \pm 14.5$	$887.8 \pm 58.6$	$976.5 \pm 3.3$	$968.9 \pm 6.8$	$819.0 \pm 142.4$	$\textbf{981.3} \pm 2.8$
cartpole	swingup	$756.9 \pm 139.8$	$729.5 \pm 185.0$	$869.9 \pm 7.0$	$794.8 \pm 113.9$	$865.3 \pm 5.2$	$827.4 \pm 20.2$	$870.0 \pm 7.1$
walker	walk	$586.1 \pm 146.6$	$808.7 \pm 67.5$	$887.8 \pm 64.4$	$\textbf{959.2} \pm 6.8$	$927.3 \pm 32.6$	$822.6 \pm 49.6$	$954.2 \pm 10.4$
reacher	easy	$526.6 \pm 224.1$	$766.1 \pm 204.1$	$916.4 \pm 45.5$	$600.0\pm203.9$	$851.8 \pm 160.0$	$942.1 \pm 27.9$	$\textbf{962.1} \pm 7.1$
finger	spin	$807.3 \pm 67.2$	$798.6 \pm 81.6$	$565.8 \pm 169.4$	$919.1 \pm 25.6$	$951.3 \pm 14.6$	$837.0 \pm 26.9$	$982.9 \pm 0.8$
quadruped	walk	$199.1 \pm 52.4$	$790.1 \pm 66.4$	$865.7 \pm 93.9$	$762.7 \pm 142.0$	$764.3 \pm 125.4$	$883.0 \pm 17.7$	$\textbf{943.3} \pm 20.4$
cartpole	swingup_sparse	$617.0 \pm 164.0$	$693.7 \pm 182.5$	$728.9 \pm 104.1$	$735.7 \pm 139.5$	$829.8 \pm 9.1$	$723.6 \pm 145.9$	$\textbf{837.5} \pm 9.3$
reacher	hard	$628.9 \pm 169.2$	$574.4 \pm 233.5$	$362.2\pm162.2$	$375.9 \pm 236.3$	$\textbf{870.5} \pm 38.1$	$853.9 \pm 170.0$	$817.5 \pm 136.7$
finger	turn_easy	$204.0 \pm 31.8$	$766.0 \pm 86.4$	$381.7 \pm 152.1$	$742.7 \pm 131.9$	$847.8 \pm 89.8$	$765.1 \pm 92.6$	$\textbf{920.8} \pm 23.9$
walker	run	$143.7\pm54.3$	$450.5 \pm 42.2$	$722.7 \pm 14.8$	$743.5 \pm 17.2$	$669.3 \pm 17.8$	$446.7 \pm 91.5$	$795.2 \pm 5.0$
cheetah	run	$69.3 \pm 104.1$	$508.0 \pm 30.6$	$522.6 \pm 99.1$	$669.8 \pm 13.4$	$662.9 \pm 19.1$	$670.3 \pm 64.1$	$\textbf{874.6} \pm 14.5$
pendulum	swingup	$332.5\pm251.6$	$\textbf{835.9} \pm 5.6$	$291.1\pm198.8$	$766.7 \pm 150.9$	$372.0 \pm 264.5$	$283.6 \pm 238.8$	$765.3 \pm 151.3$
quadruped	run	$193.9 \pm 73.4$	$457.0 \pm 11.7$	$537.0 \pm 63.3$	$486.9 \pm 29.8$	$421.8 \pm 79.0$	$414.1\pm59.0$	$\textbf{889.8} \pm 41.5$
finger	turn_hard	$108.4 \pm 16.3$	$702.2 \pm 104.5$	$211.5 \pm 143.4$	$637.6 \pm 85.1$	$707.7 \pm 159.4$	$649.8 \pm 134.5$	$879.2 \pm 20.5$
hopper	stand	$97.2 \pm 126.9$	$191.7\pm128.6$	$41.4 \pm 12.7$	$292.9 \pm 123.9$	$209.6 \pm 148.6$	$515.7 \pm 183.8$	$\pmb{843.0} \pm 82.2$
acrobot	swingup	$30.1 \pm 12.6$	$222.6 \pm 17.9$	$40.9 \pm 27.3$	$59.8 \pm 33.4$	$25.8 \pm 13.7$	$12.6 \pm 2.6$	$\textbf{245.9} \pm 37.0$
hopper	hop	$3.7\pm4.5$	$5.5\pm5.5$	$6.8\pm3.5$	$93.8 \pm 35.2$	$18.4 \pm 13.2$	$122.4 \pm 37.4$	$\textbf{238.5} \pm 28.4$

Table 14: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=1, noise dim=16.

Table 15: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=1, noise dim=64.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup cartpole cartpole walker cartpole	catch balance balance_sparse stand swingup	$\begin{array}{c} 810.6 \pm 189.0 \\ 889.0 \pm 37.5 \\ 885.8 \pm 36.0 \\ 384.4 \pm 117.1 \\ 726.1 \pm 142.4 \\ 408.2 \pm 80.8 \end{array}$	$\begin{array}{c} 965.1 \pm 2.7 \\ 967.2 \pm 11.1 \\ 938.5 \pm 35.3 \\ 923.9 \pm 29.8 \\ 853.0 \pm 5.7 \\ 828.3 \pm 71.2 \end{array}$	$903.0 \pm 78.8$ $506.0 \pm 259.4$ $741.4 \pm 120.5$ $776.7 \pm 63.6$ $874.1 \pm 9.2$ $915.6 \pm 43.6$	$\begin{array}{c} 974.7 \pm 1.4 \\ 980.5 \pm 16.9 \\ 965.1 \pm 37.6 \\ 970.7 \pm 9.7 \\ 770.9 \pm 108.4 \\ 944.7 \pm 21.5 \end{array}$	$\begin{array}{c} 853.7 \pm 176.4 \\ 978.8 \pm 18.9 \\ 875.1 \pm 117.6 \\ 965.9 \pm 4.4 \\ 855.8 \pm 3.9 \\ 940.4 \pm 7.5 \end{array}$	$\begin{array}{c} 970.6 \pm 3.1 \\ 932.7 \pm 29.9 \\ 927.2 \pm 79.9 \\ 698.4 \pm 163.0 \\ 836.5 \pm 12.7 \\ 802.7 \pm 79.2 \end{array}$	$978.2 \pm 1.0$ $999.8 \pm 0.1$ $989.9 \pm 13.1$ $980.9 \pm 3.3$ $870.7 \pm 6.0$ $964.8 \pm 2.7$
reacher finger quadruped cartpole reacher finger walker cheetah pendulum	easy spin walk swingup_sparse hard turn_easy run run swingup	$\begin{array}{c} 538.6 \pm 155.7 \\ 460.5 \pm 197.1 \\ 233.6 \pm 65.9 \\ 739.4 \pm 67.3 \\ 493.5 \pm 176.6 \\ 204.4 \pm 21.4 \\ 138.3 \pm 54.1 \\ 42.4 \pm 63.3 \\ 203.6 \pm 211.9 \end{array}$	$\begin{array}{c} 761.9 \pm 211.0 \\ 776.3 \pm 45.8 \\ 799.8 \pm 52.6 \\ 793.7 \pm 28.5 \\ 594.9 \pm 266.4 \\ 708.6 \pm 108.3 \\ 443.9 \pm 43.2 \\ 487.7 \pm 30.1 \\ 803.3 \pm 43.7 \end{array}$	$\begin{array}{c} 778.9 \pm 193.9 \\ 496.9 \pm 107.2 \\ 824.5 \pm 134.6 \\ 611.0 \pm 162.1 \\ 508.8 \pm 181.3 \\ 336.0 \pm 116.2 \\ 706.2 \pm 15.5 \\ 340.3 \pm 121.6 \\ 455.6 \pm 254.6 \end{array}$	$\begin{array}{c} 751.9 \pm 148.1 \\ 915.5 \pm 28.1 \\ 825.9 \pm 77.3 \\ 777.3 \pm 28.8 \\ 654.3 \pm 209.1 \\ 718.1 \pm 103.7 \\ 727.3 \pm 22.7 \\ 616.1 \pm 19.8 \\ \textbf{839.1} \pm 5.3 \end{array}$	$\begin{array}{c} 950.7\pm1.3\\ 950.7\pm1.3\\ 837.1\pm15.0\\ 795.8\pm154.0\\ 904.2\pm41.1\\ 548.0\pm204.3\\ 646.4\pm23.2\\ 598.4\pm19.9\\ 494.2\pm269.3 \end{array}$	$\begin{array}{c} 798.1 \pm 209.1 \\ 802.9 \pm 54.3 \\ 797.5 \pm 117.1 \\ 798.8 \pm 20.7 \\ 771.1 \pm 197.9 \\ 500.6 \pm 133.6 \\ 495.6 \pm 82.7 \\ 696.8 \pm 28.9 \\ 497.8 \pm 257.4 \end{array}$	$\begin{array}{c} \textbf{962.3} \pm 9.2\\ \textbf{964.4} \pm 17.7\\ \textbf{945.6} \pm 12.2\\ \textbf{838.2} \pm 8.9\\ \textbf{944.4} \pm 14.2\\ \textbf{929.3} \pm 25.8\\ \textbf{765.0} \pm 15.7\\ \textbf{841.1} \pm 20.8\\ \textbf{838.9} \pm 6.6 \end{array}$
quadruped finger hopper acrobot hopper	run turn_hard stand swingup hop	$\begin{array}{c} 195.3 \pm 69.0 \\ 102.0 \pm 18.7 \\ 9.6 \pm 8.7 \\ 17.3 \pm 9.1 \\ 2.6 \pm 4.3 \end{array}$	$\begin{array}{c} 438.4\pm26.7\\ 294.7\pm87.1\\ 35.3\pm28.5\\ 67.4\pm22.3\\ 1.4\pm0.9\end{array}$	$\begin{array}{c} 526.2\pm 72.2\\ 174.0\pm 67.9\\ 29.4\pm 10.5\\ 29.5\pm 24.6\\ 3.4\pm 2.4\\ \end{array}$	$\begin{array}{c} 456.8\pm 38.4\\ 597.3\pm 58.8\\ 351.2\pm 175.7\\ 94.9\pm 55.1\\ 92.9\pm 29.0\\ \end{array}$	$\begin{array}{c} 420.0\pm80.5\\ 279.0\pm159.5\\ 94.4\pm53.1\\ 13.9\pm4.4\\ 2.1\pm3.8\end{array}$	$\begin{array}{c} 419.9\pm 31.0\\ 344.3\pm 142.0\\ 580.2\pm 165.3\\ 12.7\pm 3.5\\ 120.8\pm 44.6\end{array}$	$\begin{array}{c} \textbf{863.6} \pm 39.4 \\ \textbf{908.2} \pm 16.3 \\ \textbf{903.7} \pm 19.2 \\ \textbf{275.4} \pm 32.1 \\ \textbf{245.9} \pm 14.1 \end{array}$

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
ball_in_cup	catch	$602.4 \pm 260.0$	$964.0 \pm 2.3$	$793.0 \pm 143.8$	$969.4 \pm 7.6$	$817.5 \pm 216.7$	$966.6 \pm 8.1$	$\textbf{973.7} \pm 1.5$
cartpole	balance	$868.8 \pm 55.3$	$818.9 \pm 203.4$	$768.9 \pm 222.2$	$982.9 \pm 10.4$	$926.5 \pm 47.6$	$943.7 \pm 23.0$	$\textbf{998.9} \pm 1.7$
cartpole	balance_sparse	$631.0 \pm 158.5$	$972.9 \pm 19.2$	$671.7 \pm 147.7$	$956.1 \pm 33.5$	$962.6 \pm 46.1$	$865.8 \pm 128.5$	$\textbf{990.4} \pm 11.7$
walker	stand	$315.8 \pm 91.8$	$925.4 \pm 18.7$	$863.1 \pm 72.6$	$972.1 \pm 9.2$	$953.3 \pm 18.0$	$765.7 \pm 128.3$	$980.8 \pm 2.1$
cartpole	swingup	$779.2 \pm 27.6$	$849.1 \pm 8.6$	$864.6 \pm 4.8$	$641.8 \pm 194.5$	$854.4 \pm 8.7$	$835.5 \pm 8.2$	$\textbf{874.0} \pm 5.0$
walker	walk	$239.1\pm112.6$	$879.3 \pm 46.0$	$899.3 \pm 69.0$	$954.8 \pm 6.0$	$910.8 \pm 29.5$	$754.4 \pm 158.0$	$\textbf{962.0} \pm 6.1$
reacher	easy	$297.5 \pm 157.4$	$621.3 \pm 254.4$	$704.9 \pm 288.4$	$799.6 \pm 147.2$	$935.3 \pm 31.1$	$859.9 \pm 160.6$	$\textbf{954.5} \pm 13.1$
finger	spin	$394.5 \pm 194.0$	$666.3 \pm 55.8$	$433.7 \pm 78.0$	$852.5 \pm 64.3$	$657.3 \pm 53.4$	$842.4 \pm 64.8$	$\textbf{934.7} \pm 38.2$
quadruped	walk	$209.5 \pm 49.9$	$789.1 \pm 53.4$	$803.6 \pm 117.2$	$708.7 \pm 125.6$	$567.7 \pm 222.2$	$822.5 \pm 100.1$	$\textbf{944.7} \pm 14.4$
cartpole	swingup_sparse	$665.0 \pm 104.9$	$0.6 \pm 0.7$	$212.2\pm172.6$	$517.7 \pm 225.3$	$397.2 \pm 264.0$	$650.1 \pm 194.3$	$\textbf{836.3} \pm 9.5$
reacher	hard	$712.8 \pm 152.2$	$263.2 \pm 213.3$	$242.2 \pm 177.5$	$747.7 \pm 158.7$	$693.0 \pm 223.5$	$899.3 \pm 85.7$	$\textbf{944.9} \pm 10.6$
finger	turn_easy	$210.6 \pm 46.1$	$477.5 \pm 105.4$	$237.6 \pm 62.4$	$655.3 \pm 53.6$	$459.0 \pm 118.5$	$374.4 \pm 104.9$	$\textbf{936.8} \pm 7.5$
walker	run	$26.2 \pm 2.9$	$449.9 \pm 62.1$	$703.9 \pm 15.7$	$710.6 \pm 28.9$	$643.3 \pm 14.8$	$507.0 \pm 97.3$	$766.6 \pm 12.6$
cheetah	run	$219.5 \pm 80.4$	$441.5 \pm 42.3$	$297.7 \pm 142.4$	$508.5 \pm 39.8$	$577.4 \pm 23.6$	$648.1 \pm 35.9$	$862.4 \pm 11.9$
pendulum	swingup	$85.8 \pm 78.0$	$764.4 \pm 122.2$	$741.3 \pm 140.4$	$706.9 \pm 176.5$	$349.4 \pm 271.1$	$166.2\pm197.5$	$\textbf{839.5} \pm 5.6$
quadruped	run	$181.8\pm52.9$	$446.5\pm21.8$	$508.2 \pm 61.6$	$408.6 \pm 73.9$	$315.3 \pm 72.3$	$409.5\pm50.1$	$\pmb{867.3} \pm 37.0$
finger	turn_hard	$94.3 \pm 24.0$	$213.2\pm60.2$	$126.2 \pm 21.1$	$497.7 \pm 49.4$	$278.4 \pm 110.5$	$290.4 \pm 140.9$	$\pmb{899.2} \pm 27.3$
hopper	stand	$15.1 \pm 12.4$	$28.7 \pm 23.2$	$15.2 \pm 8.7$	$261.3 \pm 153.5$	$50.6 \pm 35.8$	$524.3 \pm 197.3$	$874.0 \pm 42.3$
acrobot	swingup	$16.1 \pm 8.7$	$24.9 \pm 8.5$	$29.1 \pm 14.5$	$101.9\pm50.1$	$14.0 \pm 3.0$	$10.5 \pm 2.5$	$243.2 \pm 41.1$
hopper	hop	$0.4\pm0.2$	$0.9\pm0.1$	$1.8 \pm 2.3$	$69.5 \pm 26.3$	$0.1\pm0.2$	$119.9 \pm 28.8$	$\textbf{245.2} \pm 16.9$

Table 16: Performance of state-based DMC tasks for the compared methods in noise setting: noise std=1, noise dim=128.

Table 17: Performance of pixel-based DMC tasks for the compared methods with original clean background.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
cartpole	balance	$963.7 \pm 21.7$	$984.4 \pm 9.1$	$981.9 \pm 27.0$	$945.0 \pm 27.8$	$\textbf{990.1} \pm 5.1$	$943.2\pm67.7$	$986.8 \pm 5.8$
cartpole	balance_sparse	$617.1\pm629.2$	$984.0 \pm 29.1$	$992.7 \pm 16.2$	$995.1 \pm 6.1$	<b>997.9</b> ± 3.9	$916.3 \pm 81.8$	$972.0 \pm 47.9$
walker	stand	$923.8 \pm 63.4$	$968.3 \pm 7.7$	$\textbf{971.1} \pm 6.3$	$963.2 \pm 8.6$	$968.6 \pm 9.0$	$961.1 \pm 20.0$	$875.2 \pm 219.7$
finger	spin	$837.5 \pm 197.7$	$845.5 \pm 231.1$	$967.1 \pm 40.7$	$913.5 \pm 179.2$	$918.2 \pm 173.4$	$\textbf{973.9} \pm 29.9$	$918.6 \pm 181.1$
cartpole	swingup	$841.9 \pm 35.7$	$856.1 \pm 7.0$	$\textbf{859.5} \pm 2.5$	$815.6 \pm 21.8$	$852.3 \pm 16.0$	$856.2 \pm 15.6$	$838.6 \pm 45.7$
ball_in_cup	catch	$311.0 \pm 450.8$	$950.2 \pm 20.4$	<b>970.0</b> ± 1.9	$843.5 \pm 180.9$	$789.1 \pm 481.4$	$793.9 \pm 499.4$	$969.8 \pm 6.8$
walker	walk	$658.8 \pm 184.6$	$764.3 \pm 18.9$	$936.6 \pm 32.3$	$941.9 \pm 25.7$	$785.8 \pm 129.8$	$955.9 \pm 16.0$	$\textbf{957.2} \pm 4.1$
point_mass	easy	$\textbf{830.1} \pm 69.4$	$512.5\pm579.2$	$715.1 \pm 495.6$	$418.7 \pm 490.5$	$287.4\pm506.9$	$490.4 \pm 558.9$	$246.4 \pm 483.4$
cartpole	swingup_sparse	$732.8 \pm 18.8$	$780.7 \pm 49.3$	$758.6 \pm 42.0$	$745.0 \pm 74.7$	$766.7 \pm 45.5$	$\textbf{800.8} \pm 30.6$	$655.2 \pm 455.3$
reacher	easy	$148.5\pm107.9$	$228.3 \pm 50.2$	$\textbf{924.1} \pm 44.5$	$212.7 \pm 99.3$	$583.0 \pm 147.5$	$129.9 \pm 39.2$	$84.6 \pm 19.7$
pendulum	swingup	$19.3\pm4.7$	$814.8 \pm 26.3$	$\pmb{841.6} \pm 8.2$	$835.3 \pm 11.0$	$424.5\pm513.2$	$839.3 \pm 6.7$	$668.9 \pm 463.2$
cheetah	run	$223.9 \pm 270.5$	$583.6 \pm 98.3$	$\textbf{683.4} \pm 221.2$	$302.2\pm215.1$	$586.9 \pm 154.5$	$403.7\pm343.3$	$672.7 \pm 465.3$
walker	run	$142.9 \pm 111.1$	$480.6 \pm 32.8$	$530.1 \pm 225.6$	$\textbf{578.2} \pm 47.8$	$304.3 \pm 74.3$	$503.2 \pm 171.3$	$550.5 \pm 74.8$
hopper	hop	$18.8 \pm 21.4$	$178.1 \pm 21.2$	$108.5\pm 61.6$	$129.2\pm90.0$	$160.9 \pm 21.6$	$\textbf{203.1} \pm 52.5$	$113.3 \pm 136.8$

Table 18: Performance of pixel-based DMC tasks for the compared methods in colored images background.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
cartpole	balance	$977.8 \pm 15.4$	$975.8 \pm 14.2$	$932.5 \pm 97.6$	$933.3 \pm 28.1$	$979.3 \pm 16.8$	$982.1 \pm 12.6$	$980.5 \pm 16.8$
cartpole	balance_sparse	$781.5 \pm 484.3$	$987.1 \pm 25.3$	<b>998.2</b> ± 3.5	$825.9 \pm 480.1$	$985.1 \pm 14.6$	$998.2 \pm 4.2$	$990.8 \pm 18.9$
walker	stand	$883.3 \pm 166.4$	$938.8 \pm 57.4$	$963.4 \pm 8.1$	$958.7 \pm 9.7$	$941.3 \pm 53.9$	$949.7 \pm 47.8$	$948.4 \pm 26.4$
finger	spin	$489.8 \pm 583.7$	$981.8 \pm 2.9$	$786.8 \pm 546.2$	$670.4 \pm 504.0$	$975.3 \pm 17.4$	$902.4 \pm 207.2$	$984.4 \pm 4.0$
cartpole	swingup	$827.5 \pm 25.7$	$827.5 \pm 37.6$	$\textbf{856.6} \pm 9.8$	$799.6 \pm 11.9$	$851.0\pm7.9$	$844.1 \pm 32.5$	$707.1\pm373.0$
ball_in_cup	catch	$119.9 \pm 38.9$	$775.3 \pm 454.1$	$967.7 \pm 6.3$	$613.6 \pm 573.5$	$605.6 \pm 599.0$	$918.8 \pm 146.8$	$459.9 \pm 579.6$
walker	walk	$635.0 \pm 37.4$	$\textbf{799.8} \pm 162.9$	$907.9 \pm 19.0$	$828.2 \pm 125.8$	$701.7 \pm 158.9$	$570.0\pm620.9$	$\textbf{938.2} \pm 23.8$
point_mass	easy	$258.2 \pm 440.5$	$\textbf{710.3} \pm 493.0$	$544.9 \pm 602.0$	$527.0 \pm 406.0$	$156.3 \pm 424.4$	$524.1\pm592.5$	$670.8 \pm 468.6$
cartpole	swingup_sparse	$717.7 \pm 49.0$	$787.5 \pm 48.0$	$770.3 \pm 59.1$	$551.3 \pm 386.5$	$611.6 \pm 429.4$	$\pmb{804.3} \pm 51.0$	$797.2 \pm 54.1$
reacher	easy	$129.0 \pm 58.2$	$188.5 \pm 46.0$	$\textbf{952.0} \pm 7.6$	$164.6 \pm 73.9$	$673.9 \pm 125.1$	$195.8 \pm 41.1$	$83.0 \pm 20.9$
pendulum	swingup	$1.1\pm0.7$	$\textbf{622.7} \pm 650.5$	$414.1\pm473.9$	$174.1\pm466.3$	$70.4 \pm 128.8$	$168.7\pm452.5$	$336.1\pm570.1$
cheetah	run	$212.8 \pm 227.5$	$368.1 \pm 31.9$	$373.6 \pm 58.2$	$284.8 \pm 196.5$	$362.2\pm40.7$	$\textbf{377.3} \pm 59.0$	$207.5\pm376.6$
walker	run	$140.9 \pm 96.7$	$385.2 \pm 33.5$	$514.1 \pm 68.3$	$531.3 \pm nan$	$269.1 \pm 45.2$	$368.6 \pm 268.3$	$\textbf{543.1} \pm 43.0$
hopper	hop	$9.2\pm20.5$	$63.9 \pm 77.2$	$71.2 \pm 83.4$	$37.4 \pm 117.6$	$110.7\pm126.1$	$\boldsymbol{174.0} \pm 42.6$	$118.3\pm133.4$

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
cartpole cartpole walker finger cartpole ball_in_cup walker	balance balance_sparse stand spin swingup catch walk	$\begin{array}{c} 788.0 \pm 343.4 \\ 791.7 \pm 540.5 \\ 822.1 \pm 189.7 \\ 558.5 \pm 632.9 \\ 830.6 \pm 24.5 \\ 205.1 \pm 276.5 \\ 502.5 \pm 224.8 \end{array}$	$\begin{array}{c} 883.6 \pm 221.3 \\ 998.2 \pm 2.3 \\ 952.6 \pm 13.5 \\ 906.8 \pm 210.6 \\ 839.9 \pm 19.3 \\ 782.8 \pm 396.7 \\ 712.2 \pm 137.8 \end{array}$	$\begin{array}{c} 986.4 \pm 13.3 \\ 997.6 \pm 6.1 \\ \textbf{961.4} \pm 10.0 \\ 955.1 \pm 66.0 \\ \textbf{861.5} \pm 11.7 \\ \textbf{797.6} \pm 484.3 \\ 922.1 \pm 13.1 \end{array}$	$\begin{array}{c} 845.4\pm 262.1\\ 761.2\pm 751.8\\ 922.8\pm 84.9\\ 705.9\pm 499.9\\ 749.2\pm 115.1\\ 568.6\pm 623.2\\ 739.5\pm 495.3\\ \end{array}$	$\begin{array}{c} \textbf{986.8} \pm 6.3 \\ \textbf{999.4} \pm 0.9 \\ \textbf{955.7} \pm 23.0 \\ \textbf{972.2} \pm 16.6 \\ \textbf{855.9} \pm 12.1 \\ \textbf{739.3} \pm 483.9 \\ \textbf{640.0} \pm 46.6 \end{array}$	$\begin{array}{c} 943.3 \pm 81.8 \\ \textbf{999.4} \pm 1.0 \\ 950.5 \pm 38.0 \\ 926.4 \pm 155.9 \\ 838.9 \pm 21.3 \\ 761.1 \pm 440.6 \\ 931.5 \pm 66.0 \end{array}$	$\begin{array}{l} 985.3 \pm 6.9 \\ 985.9 \pm 16.6 \\ 952.0 \pm 20.2 \\ \textbf{982.8} \pm 4.8 \\ 507.9 \pm 409.3 \\ 585.2 \pm 580.9 \\ \textbf{938.6} \pm 26.4 \end{array}$
point_mass cartpole reacher pendulum	easy swingup_sparse easy swingup	$\begin{array}{c} 518.8 \pm 587.9 \\ 670.5 \pm 65.7 \\ 141.6 \pm 81.7 \\ 4.9 \pm 5.7 \end{array}$	$\begin{array}{c} 535.4 \pm 604.8 \\ \textbf{794.4} \pm 20.3 \\ 203.3 \pm 70.9 \\ \textbf{549.5} \pm 572.5 \end{array}$	$\begin{array}{c} \textbf{536.7} \pm 605.3 \\ \textbf{354.8} \pm 476.7 \\ \textbf{962.4} \pm 19.4 \\ \textbf{199.7} \pm 435.7 \end{array}$	$\begin{array}{c} 201.0 \pm 435.0 \\ 596.8 \pm 358.2 \\ 179.4 \pm 111.2 \\ 337.2 \pm 553.5 \end{array}$	$\begin{array}{c} 172.7 \pm 477.7 \\ 610.0 \pm 647.6 \\ 590.2 \pm 152.4 \\ 59.3 \pm 148.8 \end{array}$	$\begin{array}{c} 324.6 \pm 545.9 \\ 629.0 \pm 437.1 \\ 150.7 \pm 73.8 \\ 339.0 \pm 569.5 \end{array}$	$\begin{array}{c} 465.7 \pm 414.0 \\ 785.0 \pm 51.5 \\ 88.1 \pm 23.6 \\ 346.2 \pm 576.4 \end{array}$
cheetah walker hopper	run run hop	$\begin{array}{c} 202.0 \pm 161.9 \\ 74.7 \pm 209.4 \\ 11.5 \pm 30.9 \end{array}$	$\begin{array}{c} 369.4\pm 56.0\\ 416.2\pm 33.1\\ 120.9\pm 88.7 \end{array}$	$\begin{array}{c} \textbf{410.0} \pm 29.6 \\ \textbf{469.4} \pm 144.8 \\ \textbf{79.9} \pm 54.9 \end{array}$	$\begin{array}{c} 217.3 \pm 244.0 \\ 310.3 \pm 327.2 \\ 73.5 \pm 82.4 \end{array}$	$\begin{array}{c} 376.8 \pm 31.9 \\ 281.3 \pm 52.4 \\ \textbf{147.0} \pm 19.8 \end{array}$	$\begin{array}{c} 286.5 \pm 198.8 \\ 267.1 \pm 313.3 \\ 100.2 \pm 117.5 \end{array}$	$\begin{array}{c} 253.4 \pm 285.4 \\ 404.5 \pm 269.8 \\ 78.8 \pm 134.1 \end{array}$

Table 19: Performance of pixel-based DMC tasks for the compared methods in grayscale images background.

Table 20: Performance of pixel-based DMC tasks for the compared methods in colored video background.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
cartpole cartpole walker finger cartpole ball_in_cup walker	balance balance_sparse stand spin swingup catch walk	$\begin{array}{c} 924.7 \pm 54.5 \\ 677.1 \pm 212.8 \\ 534.5 \pm 347.5 \\ 163.9 \pm 521.6 \\ 770.6 \pm 40.8 \\ 95.7 \pm 57.6 \\ 100.6 \pm 144.0 \end{array}$	$\begin{array}{c} 956.1 \pm 45.7 \\ 986.0 \pm 22.6 \\ 943.3 \pm 17.7 \\ 978.8 \pm 12.7 \\ 850.5 \pm 7.6 \\ 885.1 \pm 119.2 \\ 845.3 \pm 83.9 \end{array}$	$\begin{array}{c} 961.8 \pm 27.4 \\ 958.2 \pm 93.6 \\ 960.0 \pm 12.1 \\ 860.1 \pm 331.9 \\ 847.6 \pm 33.3 \\ 800.5 \pm 533.7 \\ \textbf{921.6} \pm 30.3 \end{array}$	$\begin{array}{c} 760.6 \pm 386.2 \\ 990.7 \pm 14.0 \\ 951.0 \pm 5.3 \\ 905.2 \pm 177.5 \\ 301.0 \pm 347.0 \\ 766.6 \pm 412.8 \\ 651.9 \pm 458.0 \end{array}$	$\begin{array}{c} 979.4 \pm 8.7 \\ 986.8 \pm 13.7 \\ \textbf{961.8} \pm 10.8 \\ 962.1 \pm 27.0 \\ 847.0 \pm 11.9 \\ \textbf{935.9} \pm 55.0 \\ 684.6 \pm 151.0 \end{array}$	$\begin{array}{c} \textbf{985.6} \pm 22.1 \\ \textbf{997.5} \pm 6.4 \\ \textbf{916.2} \pm 128.5 \\ \textbf{982.2} \pm 4.6 \\ \textbf{861.8} \pm 33.1 \\ \textbf{365.6} \pm 262.3 \\ \textbf{667.6} \pm 452.6 \end{array}$	$\begin{array}{c} 984.1 \pm 8.9 \\ 993.6 \pm 20.2 \\ 957.7 \pm 5.6 \\ 908.8 \pm 212.0 \\ 832.9 \pm 72.4 \\ 885.4 \pm 94.9 \\ 912.6 \pm 26.2 \end{array}$
point_mass cartpole reacher pendulum	easy swingup_sparse easy swingup	$\begin{array}{c} 2.2 \pm 3.8 \\ 0.0 \pm 0.0 \\ 203.8 \pm 84.7 \\ 6.3 \pm 8.1 \end{array}$	$\begin{array}{c} \textbf{889.8} \pm 5.7 \\ \textbf{731.2} \pm 91.1 \\ \textbf{188.6} \pm 26.2 \\ \textbf{733.2} \pm 351.4 \end{array}$	$\begin{array}{c} 25.2 \pm 31.6 \\ 10.3 \pm 14.5 \\ \textbf{949.2} \pm 45.9 \\ 452.2 \pm 453.2 \end{array}$	$\begin{array}{c} 281.4 \pm 670.6 \\ 92.1 \pm 292.3 \\ 175.2 \pm 76.3 \\ 334.7 \pm 549.6 \end{array}$	$\begin{array}{l} 889.6 \pm 10.1 \\ \textbf{734.0} \pm 97.4 \\ 249.6 \pm 66.1 \\ 85.3 \pm 93.9 \end{array}$	$\begin{array}{c} 365.6 \pm 592.9 \\ 573.0 \pm 429.4 \\ 183.3 \pm 59.1 \\ 519.2 \pm 419.7 \end{array}$	$\begin{array}{c} 178.3 \pm 477.2 \\ 675.7 \pm 240.7 \\ 86.4 \pm 18.0 \\ 654.8 \pm 447.6 \end{array}$
cheetah walker hopper	run run hop	$\begin{array}{c} 53.0 \pm 59.8 \\ 54.9 \pm 23.3 \\ 0.2 \pm 0.4 \end{array}$	$\begin{array}{c} 341.2 \pm 34.0 \\ 367.0 \pm 33.4 \\ 108.2 \pm 80.7 \end{array}$	$\begin{array}{c} \textbf{395.4} \pm 21.2 \\ \textbf{492.7} \pm 135.4 \\ 41.2 \pm 45.3 \end{array}$	$\begin{array}{c} 220.3 \pm 49.8 \\ 148.5 \pm 153.7 \\ 27.6 \pm 46.0 \end{array}$	$\begin{array}{l} 348.8 \pm 8.4 \\ 244.1 \pm 70.8 \\ \textbf{136.1} \pm 18.9 \end{array}$	$\begin{array}{c} 306.7 \pm 37.8 \\ 115.8 \pm 103.8 \\ 23.0 \pm 37.0 \end{array}$	$\begin{array}{c} 315.7 \pm 21.0 \\ 168.8 \pm 170.5 \\ 27.5 \pm 50.0 \end{array}$

Table 21: Performance of pixel-based DMC tasks for the compared methods in grayscale video background.

Task		DBC	MICo	RAP	RDBC	SAC	DeepMDP	SimSR
cartpole cartpole walker finger cartpole ball_in_cup	balance balance_sparse stand spin swingup catch	$\begin{array}{c} 951.9\pm 36.1\\ 875.0\pm 134.2\\ 643.8\pm 253.2\\ 0.0\pm 0.0\\ 849.1\pm 17.5\\ 118.0\pm 39.2\\ 146.5\pm 84.5\\ \end{array}$	$\begin{array}{c} 977.2 \pm 6.6 \\ 996.3 \pm 2.2 \\ 947.8 \pm 25.7 \\ 979.9 \pm 7.7 \\ 854.3 \pm 11.1 \\ 930.5 \pm 28.0 \\ 763.3 \pm 258.4 \end{array}$	$\begin{array}{c} 983.2 \pm 12.6 \\ 956.8 \pm 67.6 \\ 971.0 \pm 3.9 \\ 971.4 \pm 21.4 \\ 861.7 \pm 8.0 \\ 958.9 \pm 23.3 \\ 913.6 \pm 22.1 \end{array}$	$\begin{array}{c} 882.2 \pm 271.8 \\ 803.4 \pm 545.8 \\ 961.1 \pm 3.1 \\ 974.8 \pm 13.1 \\ 389.8 \pm 388.4 \\ 865.0 \pm 144.5 \\ 885.0 \pm 46.9 \end{array}$	$\begin{array}{c} 973.3 \pm 17.5 \\ 991.0 \pm 15.2 \\ 961.5 \pm 3.7 \\ 976.9 \pm 11.8 \\ 858.0 \pm 14.2 \\ \textbf{962.0} \pm 3.6 \\ 650.4 \pm 101.2 \end{array}$	$\begin{array}{c} \textbf{990.7} \pm 6.7 \\ \textbf{996.4} \pm 8.9 \\ \textbf{946.6} \pm 64.9 \\ \textbf{982.7} \pm 2.9 \\ \textbf{863.4} \pm 12.3 \\ 6\textbf{93.5} \pm 259.7 \\ \textbf{869.2} \pm 255.7 \end{array}$	$\begin{array}{c} 962.9 \pm 57.5\\ 990.4 \pm 23.0\\ 957.0 \pm 8.7\\ 902.6 \pm 261.3\\ 854.7 \pm 51.0\\ 907.4 \pm 70.0\\ 907.8 \pm 48.7 \end{array}$
point_mass cartpole reacher pendulum	easy swingup_sparse easy swingup	$\begin{array}{c} 140.3 \pm 84.3 \\ \hline 156.8 \pm 428.7 \\ 0.0 \pm 0.0 \\ 162.8 \pm 115.3 \\ 33.9 \pm 90.5 \end{array}$	$\begin{array}{c} \textbf{703.3} \pm 238.4 \\ \textbf{893.0} \pm 11.0 \\ \textbf{803.0} \pm 26.9 \\ 171.1 \pm 69.7 \\ \textbf{695.5} \pm 391.8 \end{array}$	$913.0 \pm 35.1$ $439.9 \pm 356.9$ $387.9 \pm 440.7$ $968.2 \pm 14.7$ $708.2 \pm 290.2$	$\begin{array}{c} 70.8 \pm 195.6 \\ 339.2 \pm 400.9 \\ 166.9 \pm 42.4 \\ 269.1 \pm 596.6 \end{array}$	$\begin{array}{c} 885.1 \pm 28.4 \\ 643.4 \pm 447.7 \\ 268.3 \pm 83.7 \\ 207.5 \pm 459.2 \end{array}$	$\begin{array}{c} 809.2 \pm 23.3 \\ \hline 709.9 \pm 492.9 \\ 643.8 \pm 449.4 \\ 209.2 \pm 53.8 \\ \hline \textbf{836.2} \pm 17.7 \end{array}$	$\begin{array}{c} 541.6 \pm 596.3 \\ 490.0 \pm 555.5 \\ 84.9 \pm 16.5 \\ 835.6 \pm 15.4 \end{array}$
cheetah walker hopper	run run hop	$\begin{array}{c} 77.4 \pm 62.4 \\ 63.1 \pm 17.1 \\ 0.1 \pm 0.1 \end{array}$	$\begin{array}{c} 366.9 \pm 12.9 \\ 393.4 \pm 16.2 \\ 117.3 \pm 81.9 \end{array}$	$\begin{array}{l} \textbf{404.8} \pm 20.8 \\ \textbf{544.0} \pm 26.3 \\ 77.3 \pm 53.2 \end{array}$	$\begin{array}{c} 170.2 \pm 194.0 \\ 180.3 \pm 109.5 \\ 38.2 \pm 66.9 \end{array}$	$\begin{array}{c} 359.9 \pm 25.5 \\ 244.8 \pm 45.7 \\ \textbf{121.1} \pm 84.2 \end{array}$	$\begin{array}{c} 343.2 \pm 69.7 \\ 175.8 \pm 106.2 \\ 18.3 \pm 49.4 \end{array}$	$\begin{array}{c} 343.4 \pm 18.5 \\ 270.5 \pm 190.6 \\ 0.9 \pm 0.9 \end{array}$



Figure 10: Individual state-based task performance.



Figure 11: Individual state-based task performance.



Figure 12: Individual state-based task performance.



Figure 13: Individual state-based task performance.



Figure 14: Individual state-based task performance.



Figure 15: Individual state-based task performance.



Figure 16: Individual state-based task performance.



Figure 17: Individual state-based task performance.



Figure 18: Individual state-based task performance.



Figure 19: Individual pixel-based task performance.



#### Noise setting: Natural Image (Grayscale)

Figure 20: Individual pixel-based task performance.



#### Noise setting: Natural Images

Figure 21: Individual pixel-based task performance.



#### Noise setting: Natural Video (Grayscale)

Figure 22: Individual pixel-based task performance.



Figure 23: Individual pixel-based task performance.











Figure 26: Denoising factor (ID evaluation setting) on isolated metric encoder  $\tilde{\phi}$ .



Figure 27: Denoising factor (OOD evaluation setting) on isolated metric encoder  $\tilde{\phi}$ .



Figure 28: Denoising factor (ID evaluation setting) on isolated metric encoder  $\tilde{\phi}$ .



Figure 29: Denoising factor (OOD evaluation setting) on isolated metric encoder  $\tilde{\phi}$ .