
Learning Latent Dynamic Robust Representations for World Models

Anonymous Authors¹

Abstract

Visual Model-Based Reinforcement Learning (MBRL) promises to encapsulate agent's knowledge about the underlying dynamics of the environment, enabling learning a world model as a useful planner. However, top MBRL agents such as Dreamer often struggle with visual pixel-based inputs in the presence of exogenous or irrelevant noise in the observation space, due to failure to capture task-specific features while filtering out irrelevant spatio-temporal details. To tackle this problem, we apply a spatio-temporal masking strategy, a bisimulation principle, combined with latent reconstruction, to capture endogenous task-specific aspects of the environment for world models, effectively eliminating non-essential information. Joint training of representations, dynamics, and policy often leads to instabilities. To further address this issue, we develop a Hybrid Recurrent State-Space Model (HRSSM) structure, enhancing state representation robustness for effective policy learning. Our empirical evaluation demonstrates significant performance improvements over existing methods in a range of visually complex control tasks such as Maniskill (Gu et al., 2023) with exogenous distractors from the Matterport environment.

1. Introduction

Model-Based Reinforcement Learning (MBRL) utilizes predictive models to capture endogenous dynamics of the world, to be able to simulate and forecast future scenarios, enhancing the agent's decision making abilities by leveraging imagination and prediction in visual pixel-based contexts (Hafner et al., 2019a; Kalweit & Boedecker, 2017; Hafner et al., 2019b; Ha & Schmidhuber, 2018; Janner et al., 2020). Most importantly, these world models such as recurrent state-space model (RSSM) (Hafner et al., 2019b), enable

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

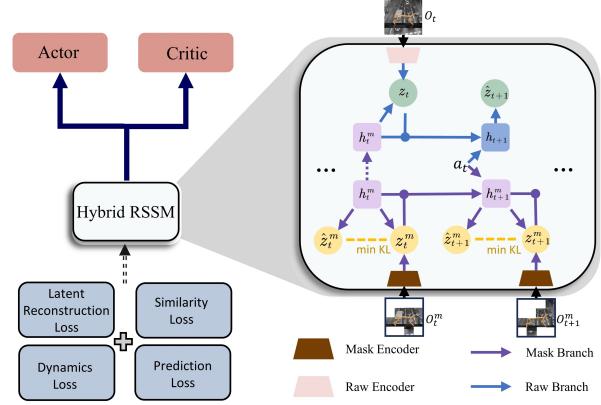


Figure 1. Our framework is composed of a Hybrid-RSSM and actor-critic architecture. Hybrid-RSSM learns robust representations and dynamics through four distinct objectives: latent reconstruction, which aligns features between masked and raw observations; similarity loss based on the bisimulation principle; and two additional objectives same as in Dreamer series (Hafner et al., 2020; 2023).

agents to understand and represent dynamics in the learnt representation space, consisting of task specific information with the hope to have filtered out exogenous or irrelevant aspects from the observations, leading to superior performance compared to model-free RL algorithms. However, most MBRL methods face challenges in environments with large amounts of unpredictable or irrelevant exogenous observations (Burda et al., 2018; Efroni et al., 2022; 2021).

Arguably, the Dreamer series of algorithms (Hafner et al., 2019a; 2020; 2023) are probably the most effective and representative class of MBRL approaches where agents learn representations and dynamics in latent space by minimizing reconstruction errors. Most MBRL approaches such as Dreamer often includes a forward dynamics model to predict observations and a reward model that evaluates the potential of future states. Recent works however have shown the ineffectiveness of forward dynamics based models when learning from exogenous observation based visual inputs (Efroni et al., 2022; Lamb et al., 2022; Islam et al., 2022). This is because in noisy environments, emphasis on reconstruction can lead to disproportionate focus on irrelevant details such as textures or noise, at the expense of smaller but task-relevant elements. This can result in inaccuracies in the dynamics model (Xiao et al., 2019; Asadi et al., 2019)

and overfitting to specific environmental traits (Zhang et al., 2020a), leading to compounded errors in latent space world models for planning.

While a body of work has addressed exogenous noise, primarily in reward-free (Efroni et al., 2021; 2022; Lamb et al., 2022) or offline visual settings for model-free RL (Islam et al., 2022), only limited research has explored model-based agents in the context of exogenous noise. These studies have developed in a way of decoder-free matter, *i.e.*, excluding pixel-level reconstruction, to mitigate reconstruction issues, but they still face significant challenges: either lacking in capturing task-specific information (Deng et al., 2022; Okada & Taniguchi, 2021), not being robust against various noise types (Fu et al., 2021), or sensitive to hyperparameters (Zhu et al., 2023; Henderson et al., 2018). This work is therefore primarily driven the question :

How to learn sufficiently expressive state representation for a world model without the reliance of the pixel-level reconstruction?

In principle, the ideal representation objective for model-based planners should address two desired criterion : i) effectively capturing task-relevant endogenous dynamics information, and ii) be robust and compact enough to filter exogenous task irrelevant details. Despite several prior works trying to address this (Lamb et al., 2022; Efroni et al., 2021; Islam et al., 2022) in reward free settings, these works do not show effectiveness of the learnt representation for use in world models. We address this question through the promising approach of bisimulation principle (Ferns et al., 2011; Castro, 2020; Zhang et al., 2020b; Castro et al., 2021; Zang et al., 2022a), learning representations specific to task objectives that can reflect state behavioral similarities. However, the effectiveness of the bisimulation metric heavily relies on the accuracy of the dynamics model (Kemertas & Aumentado-Armstrong, 2021). Under an approximate dynamics model, the state representation guided by the bisimulation principle may be task-specific but not necessarily compact, indicating a gap in the bisimulation principle’s ability to foster expressive state representations for robust model-based agents.

To effectively apply bisimulation principle in world models, we propose to develop a new architecture - the Hybrid-RSSM (HRSSM). This architecture employs a masking strategy to foster more compact latent representations, specifically targeting the integration of the bisimulation principle to improve the efficiency and effectiveness of the model. Our Hybrid-RSSM consists of two branches: 1) the raw branch, which processes original interaction sequences, and 2) the mask branch, which handles sequences that have been transformed using a masking strategy. This masking, involving cubic sampling of observation sequences, is designed to reduce spatio-temporal redundancy in natural signals. A key

feature of our approach is the reconstruction of masked observations to match the latent features from the raw branch in the latent representation space, not in pixel space. This ensures semantic alignment for both branches. Meanwhile, we incorporate a similarity-based objective, in line with the bisimulation principle, to integrate differences in immediate rewards and dynamics into the state representations.

Furthermore, to enhance training stability and minimize potential representation drift, the raw and mask branches share a unified historical information representation. This holistic structure defines our Hybrid Recurrent State Space Model (HRSSM), serving as a world model that leverages the strengths of the RSSM architecture to effectively capture task-specific information, guided by the bisimulation principle, and efficiently condense features through mask-based latent reconstruction. Our primary contributions are summarized as follows.

- We introduce Hybrid RSSM that integrates masking-based latent reconstruction and the bisimulation principle into a model-based RL framework, enabling the learning of task-relevant representations capturing endogenous dynamics.
- We study the roles of masking-based latent reconstruction and the bisimulation principle in model-based RL with empirical and theoretical analysis.
- Empirically, we evaluate our Hybrid-RSSM and actor-critic architecture by integrating it into the DreamerV3 framework, and show that the resulting model can be used to solve complex tasks consisting of a variety of exogenous visual information.

2. Related Work

MBRL and World Model Model-based Reinforcement Learning (MBRL) stands as a prominent subfield in Reinforcement Learning, aiming to optimize total reward through action sequences derived from dynamics and reward models (Sutton, 1990; Hamrick, 2019). Early approaches in MBRL typically focus on low-dimensional and compact state spaces (Williams et al., 2017; Janner et al., 2019; 2020), yet they demonstrated limited adaptability to more complex high-dimensional spaces. Recent efforts (Hafner et al., 2019b;a; 2020; 2023; Hansen et al., 2022a; Rafailov et al., 2021; Gelada et al., 2019) have shifted towards learning world models for these intricate spaces, utilizing visual inputs and other signals like scalar rewards. These methods enable agents to simulate behaviors in a conceptual model, thereby reducing the reliance on physical environment interactions. As a notable example, Dreamer (Hafner et al., 2019a; 2020; 2023) learns recurrent state-space models (RSSM) and the latent state space via reconstruction

losses, though achieving a good performance in conventional environments yet fails in environments with much exogenous noise.

Model-based Representation Learning Many recent MBRL methods start to integrate state representation learning into their framework to improve the robustness and efficiency of the model. Some approaches formulations rely on strong assumptions (Gelada et al., 2019; Agarwal et al., 2020). Some approaches learn world model via requiring latent temporal consistency (Zhao et al., 2023; Hansen et al., 2022b; 2023). Some approaches develop upon Dreamer architecture, combining the transformer-based masked auto-encoder (Seo et al., 2023a), extending Dreamer by explicitly modeling two independent latent MDPs that represent useful signal and noise, respectively (Fu et al., 2021; Wang et al., 2022a), optimizing the world model by utilizing mutual information (Zhu et al., 2023), regularizing world model via contrastive learning (Okada & Taniguchi, 2021; Poudel et al., 2023) and prototype-based representation learning (Deng et al., 2022). Unlike other approaches that either neglect reward significance or are limited by modeling predefined noise form, our approach learns robust representations and dynamics effectively by incorporating reward-aware information and masking strategy, we provide a more detailed comparison and additional related works in Appendix C.

3. Preliminaries

MDP The standard Markov decision process (MDP) framework is given by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, with state space \mathcal{S} , action space \mathcal{A} , reward function $r(s, a)$ bounded by $[R_{\min}, R_{\max}]$, a discount factor $\gamma \in [0, 1]$, and a transition function $P(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow \Delta \mathcal{S}$ that decides the next state, where the transition function can be either deterministic, i.e., $s' = P(s, a)$, or stochastic, i.e. $s' \sim P(\cdot | s, a)$. In the sequel, we use P_s^a to denote $P(\cdot | s, a)$ or $P(s, a)$ for simplicity. The agent in the state $s \in \mathcal{S}$ selects an action $a \in \mathcal{A}$ according to its policy, mapping states to a probability distribution on actions: $a \sim \pi(\cdot | s)$. We make use of the state value function $V^\pi(s) = \mathbb{E}_{\mathcal{M}, \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s]$ to describe the long term discounted reward of policy π starting at the state s , where $\mathbb{E}_{\mathcal{M}, \pi}$ denotes expectations under $s_0 \sim P_0$, $a_t \sim \pi(\cdot | s_t)$, and $s_{t+1} \sim P_s^a$. And the goal is to learn a policy π that maximizes the sum of expected returns $\mathbb{E}_{\mathcal{M}, \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s]$.

Visual RL and Exogenous noise We address visual reinforcement learning (RL) where the agent perceives high-dimensional pixel images as observations, represented by $o_t \sim P(o_t | o_{<t}, a_{<t})$. These observations are mapped into a lower-dimensional space via a transformation \mathcal{T} and an encoder \mathcal{E} , i.e., $\mathcal{T} \circ \mathcal{E} : \mathcal{O} \rightarrow \mathcal{X}$, then generating a latent state in a latent space: $\zeta_t \in \mathcal{Z}$ through a world model. The agent's

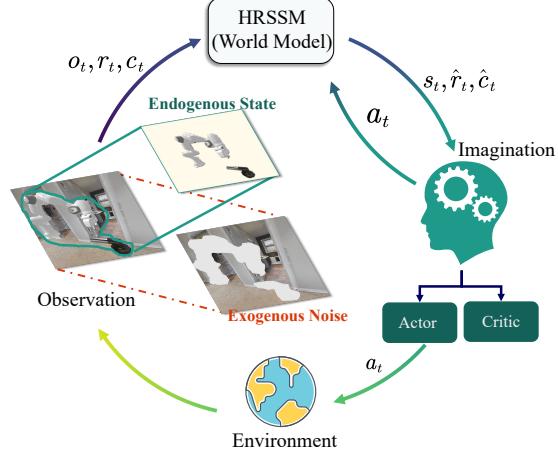


Figure 2. The entire pipeline of our framework in the presence of exogenous information. The HRSSM processes the observations into the latent space, enabling the agent to learn control within this space. Subsequently, the policy network generates actions for interacting with the environment.

actions follow a policy distribution $\pi(a | \zeta)$ under this latent state space. We introduce a setting with exogenous noise, where observations come from a mix of controllable endogenous states $s_t \in \mathcal{S}$ and uncontrollable exogenous noise $\xi_t \in \Xi$. Here, ζ_t is composed of these two components, with transitions $P(\zeta_t | \zeta_{<t}, a_{<t}) = P(s_t | s_{<t}, a_{<t})P(\xi_t | \xi_{<t})$, and rewards $r(\zeta_t, a_t) = r(s_t, a_t)$. We strive to compress latent state ζ_t by maximizing endogenous state s_t and minimizing exogenous noise ξ_t , deriving an “exogenous-free” policy, essentially $\pi(a | \zeta) \approx \pi(a | s)$. Under a mild assumption of existing mapping function ϕ_\star from the observation $o \in \mathcal{O}$ to the endogenous state $s \in \mathcal{S}$, for any o_1 and o_2 , if $\phi_\star(o_1) = \phi_\star(o_2)$, then $\pi(\cdot | o_1) = \pi(\cdot | o_2)$. The primary goal is to learn a world model that can discard exogenous noise and learn exo-free policy to improve the sample efficiency and robustness.

4. Method

In this section, we describe our overall approach of integrating the masking strategy and bisimulation principle in model-based RL methods, to learn effective world models for planning. We show that our method can be adapted to learn effective representations in the presence of exogenous noise, and the resulting planner can be used to solve complex tasks, building on the DreamerV3 (Hafner et al., 2023). The whole pipeline is shown in Figure 2.

Modified Dreamer Architecture Dreamer utilizes a recurrent state space model (RSSM) (Hafner et al., 2019b) for differentiable dynamics, learning representations of sensory inputs through backpropagation of Bellman errors from imagined trajectories. Its training process involves: optimizing the RSSM, training a policy using latent imaginations,

165 and applying this policy in the real environment. This cycle
 166 repeats until the desired policy performance is achieved.
 167 The RSSM includes several crucial components:

$$\begin{aligned}
 169 \text{Recurrent model: } & h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1}) \\
 170 \text{Representation model: } & z_t \sim q_\phi(z_t | h_t, o_t) \\
 171 \text{Transition predictor: } & \hat{z}_t \sim p_\phi(\hat{z}_t | h_t) \\
 172 \text{Reward predictor: } & \hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t) \\
 173 \text{Continue predictor: } & \hat{c}_t \sim p_\phi(\hat{c}_t | h_t, z_t) \\
 174 \text{Decoder: } & \hat{o}_t \sim p_\phi(\hat{o}_t | h_t, z_t), \\
 175 & \quad (1)
 \end{aligned}$$

176 where o_t is the sensory input, z_t the stochastic representation,
 177 h_t the recurrent state, \hat{o}_t the reconstructed input, and
 178 \hat{r}_t and \hat{c}_t are the predicted reward and the episode continuation
 179 flag. While the decoder network is crucial in Dreamer
 180 for learning environment dynamics, its reliance on recon-
 181 structing high-dimensional sensory inputs like pixels causes
 182 computational inefficiency, which arises from recovering
 183 unnecessary, control-irrelevant visual elements such as back-
 184 ground noise, impeding policy learning in environments
 185 with distractions. Prior works have explored how to recover
 186 the full endogenous latent states, by ignoring exogenous
 187 noise (Islam et al., 2022); however, effectively recovering
 188 endogenous dynamics for model-based planning remains
 189 unaddressed. We aim to develop a method for recovering
 190 these dynamics for model-based planning. Simply omitting
 191 pixel reconstruction from Dreamer, as suggested by (Hafner
 192 et al., 2019a), results in inadequate performance. Therefore,
 193 we propose modifying Dreamer to preserve accurate dynamics
 194 and enhance its awareness of essential downstream task
 195 features, while reducing dependency on reconstruction.
 196

197 4.1. Learning latent representation and dynamics

198 In visual control tasks, our state representation concentrates
 199 on two key aspects: (i) visual inputs includes much spatio-
 200 temporal redundancy, and (ii) the encapsulation of behav-
 201 iorally relevant information for the task. We introduce two
 202 novel components: masking-based latent reconstruction and
 203 similarity-based representation. The former filters out redun-
 204 dant spatiotemporal data while preserving semantic useful
 205 environmental knowledge. The latter, aligning with the bisimulation
 206 principle, retains task-specific information within the world model.
 207 This approach results in latent representations that are concise and effective.

208 Notably, our method may not recover the full endogenous
 209 dynamics, but can still be exo-free, distinguishing from
 210 other works (Lamb et al., 2022). Our key contribution is
 211 demonstrating adaptability to MBRL methods for planning,
 212 an area not fully addressed by prior research. We include
 213 detailed analysis of our proposed methodology in section
 214 5. To keep the notation succinct, we will replace ζ with s
 215 since our goal is to disregard ξ and we will ensure to remind
 216

217 readers of this when necessary.

Masking strategy Our goal is to design world models
 218 for planning that can be effective in the presence of visual
 219 exogenous information. To do this, we employ a masking
 220 strategy to reduce the spatio-temporal redundancy for
 221 enhanced control task representations. In visual RL tasks,
 222 previous works (Tong et al., 2022; Wei et al., 2022) indicate
 223 that significant spatio-temporal redundancy can be removed
 224 via masking based reconstruction methods. Consequently,
 225 we randomly mask a portion of pixels in the input observation
 226 sequence across its spatial and temporal dimensions. For a series
 227 of K environmental interaction samples $\{o_t, a_t, r_t\}_{t=1}^K$, we transform the observation
 228 sequence $\mathbf{o} = \{o_t\}_{t=1}^K \in \mathbb{R}^{K \times H \times W \times C}$ into cuboid patches
 229 $\hat{\mathbf{o}} = \{\hat{o}_t\}_{t=1}^K \in \mathbb{R}^{k P_K \times h P_H \times w P_W \times C}$, where the patch
 230 size is $(P_K \times P_H \times P_W)$ and $k = K/P_K$, $h = H/P_H$,
 231 $w = W/P_W$ are the number of patches along each dimension.
 232 We then randomly mask a fraction m of these cuboid
 233 patches to capture the most essential spatio-temporal infor-
 234 mation while discarding spatio-temporal redundancies.
 235 Subsequently, both the masked and original sequences are
 236 encoded to latent encoding space using an encoder and a
 237 momentum encoder respectively, where the momentum
 238 encoder is updated using an exponential moving average
 239 (EMA) from the masked sequence’s encoder.

Behavioral update operator To capture the task relevant
 240 information for control tasks, we adopt a similarity-based
 241 objective following the bisimulation principle (Ferns et al.,
 242 2012b;a), which requires the learnt representation to be
 243 aware of the reward and dynamics similarity between states.
 244 Our mask-based behavioral update operator, for masked and
 245 original sequences can be written as :

$$\mathcal{F}^\pi d(s_i, s_j^m) = |r_{s_i}^\pi - r_{s_j}^\pi| + \gamma \mathbb{E}_{\substack{s_{i+1} \sim \hat{P}_{s_i}^\pi, \\ s_{j+1}^m \sim \hat{P}_{s_j}^\pi}} [d(s_{i+1}, s_{j+1}^m)], \quad (2)$$

246 where s_j^m and s_i represent latent states of the mask branch
 247 and the raw branch, respectively, with $\hat{P}_{s_j}^\pi$ and $\hat{P}_{s_i}^\pi$ denoting
 248 their approximated latent dynamics, and d is the cosine
 249 distance to measure the difference between latent states.

We use equation 2 to minimize bisimulation error for learning representation. However, this process involves sampling from latent dynamics, which, when coupled with the simultaneous learning of representations, dynamics, and policies in the world model, can lead to instabilities that adversely impact dynamics learning and consequently, bisimulation training. Therefore, we develop a hybrid RSSM specifically to address complex tasks, providing a level of stability in MBRL methods, which otherwise is typically difficult due to the complexities associated with training joint objectives.

Hybrid RSSM We first follow the conventional setting of RSSM in DreamerV3 to build in the masked encoding space,

Table 1. Model components of our hybrid structure. EMA means the corresponding model is updated via exponential moving average. Gradient back-propagates through mask models and reward/continue predictor.

| | |
|---|---|
| Mask Encoder: $e_t^m = \mathcal{E}_\phi(o_t^m)$ | EMA Encoder: $e_t = \mathcal{E}'_\phi(o_t)$ |
| Mask Posterior model: $z_t^m \sim q_\phi(z_t^m h_t^m, e_t^m)$ | EMA Posterior model: $z_t \sim q'_\phi(z_t h_t^m, e_t)$ |
| Mask Recurrent model: $h_t^m = f_\phi(h_{t-1}^m, z_{t-1}^m, a_{t-1})$ | EMA Recurrent model: $h_t = f'_\phi(h_{t-1}^m, z_{t-1}, a_{t-1})$ |
| Mask Transition predictor: $\hat{z}_t^m \sim p_\phi(\hat{z}_t^m h_t^m)$ | EMA Transition predictor: $\hat{z}_t \sim p'_\phi(\hat{z}_t h_t)$ |
| Reward predictor: $\hat{r}_t \sim p_\phi(\hat{r}_t h_t^m, z_t^m)$ | Continue predictor: $\hat{c}_t \sim p_\phi(\hat{c}_t h_t^m, z_t^m)$ |

i.e., a mask encoder $e_t^m = \mathcal{E}_\phi(o_t^m)$ to encode the masked observation, a mask posterior model $z_t^m \sim q_\phi(z_t^m | h_t^m, e_t^m)$ and a mask recurrent model $h_t^m = f_\phi(h_{t-1}^m, z_{t-1}^m, a_{t-1})$ to incorporate temporal information into representations, and a mask transition predictor $\hat{z}_t^m \sim p_\phi(\hat{z}_t^m | h_t^m)$ to model the latent dynamics, where the concatenation of the mask recurrent state h_t^m and the mask posterior state z_t^m forms the mask latent state $s_t^m := [h_t^m; z_t^m]$. We train the dynamics model by minimizing the KL divergence between the posterior state z_t^m and the predicted prior state \hat{z}_t^m , and employ free bits (Kingma et al., 2016; Hafner et al., 2023), formulated as:

$$\begin{aligned} \mathcal{L}_{\text{dyn}}(\phi) &:= \beta_1 \max(1, \mathcal{L}_1(\phi)) + \beta_2 \max(1, \mathcal{L}_2(\phi)) \\ \mathcal{L}_1(\phi) &:= \text{KL}\left[\text{sg}(q_\phi(z_t^m | h_t^m, e_t^m)) \| p_\phi(\hat{z}_t^m | h_t^m)\right] \\ \mathcal{L}_2(\phi) &:= \text{KL}\left[q_\phi(z_t^m | h_t^m, e_t^m) \| \text{sg}(p_\phi(\hat{z}_t^m | h_t^m))\right] \end{aligned} \quad (3)$$

where `sg` means stopping gradient, and the values of β_1 and β_2 are set to 0.5 and 0.1, respectively, following the default configuration in DreamerV3. For now, we only construct the network of the masked sequence, but without the utilization of the original sequence. If the raw branch utilizes a different RSSM structure from the mask one, merging these complex networks could lead to training instability and representation drift. To address this, we require the raw branch and the mask branch share the same historical representation, ensuring alignment between both branches for temporal prediction. Therefore, for the raw branch, we conduct the posterior state as $z_t \sim q'_\phi(z_t | h_t^m, e_t)$, the recurrent state $h_t = f'_\phi(h_{t-1}^m, z_{t-1}, a_{t-1})$ with the historical representation from the mask branch, and the prior state $\hat{z}_t \sim p'_\phi(\hat{z}_t | h_t)$. Additionally, we define the latent state of raw branch as $s_t := [h_t^m; z_t]$ and the sampled latent state of RSSM as $\hat{s}_t = [h_t; \hat{z}_t]$. The networks q'_ϕ , f'_ϕ , and p'_ϕ are all updated using EMA from the mask branch.

We use latent reconstruction to align the feature between the masked and original ones, to disregard the unnecessary spatiotemporal redundancies, following the research within the field of computer vision (He et al., 2022; Feichtenhofer et al., 2022) that considering high-dimensional image space consists tremendous spatiotemporal redundancies. We apply a linear projection and ℓ_2 -normalize the latent state s_t and \hat{s}_t^m to obtain \bar{s}_t and \bar{s}_t^m respectively to ensure numerical stability and then compute the reconstruction loss, which

can be formulated as:

$$\mathcal{L}_{\text{rec}}(\phi) := \text{MSE}(\bar{s}_t, \bar{s}_t^m). \quad (4)$$

Meanwhile, we can minimize the bisimulation error and formulate the similarity loss to capture the task-relevant information as:

$$\begin{aligned} \mathcal{L}_{\text{sim}} &:= (d(s_i, s_j^m) - \mathcal{F}^\pi d(s_i, s_j^m))^2 \\ &= \left(d(s_i, s_j^m) - \left(|r_{s_i}^\pi - r_{s_j}^\pi| + \gamma d(\hat{s}_{i+1}, \hat{s}_{j+1}^m)\right)\right)^2, \end{aligned} \quad (5)$$

where d is the cosine distance, \hat{s}_{i+1} and \hat{s}_{j+1}^m are sampled from RSSMs.

Reward Prediction and Continue Prediction Following DreamerV3 (Hafner et al., 2023), we train the reward predictor via the symlog loss and the continue predictor via binary classification loss, to predict the reward and the episode is termination or not, they compose the prediction loss as:

$$\mathcal{L}_{\text{pred}}(\phi) := -\ln p_\phi(r_t | s_t^m) - \ln p_\phi(c_t | s_t^m). \quad (6)$$

Gradient backpropagation occurs exclusively through the mask branch, updating the representation. Consequently, we utilize only the masked latent state s_t^m for predicting both terms. Unlike the methods in Equations 4 and 5, we employ un-normalized features for prediction. Empirically, this approach enhances the model's stability and sample-efficiency, as detailed in Appendix D.

Overall The main components of our hybrid structure are illustrated in Table 1. The total loss is:

$$\mathcal{L}(\phi) := \mathbb{E}_{q_\phi} \left[\sum_{t=1}^T (\mathcal{L}_{\text{dyn}}(\phi) + \mathcal{L}_{\text{rec}}(\phi) + \mathcal{L}_{\text{sim}}(\phi) + \mathcal{L}_{\text{pred}}(\phi)) \right], \quad (7)$$

All components are optimized concurrently, with the joint minimization of the loss function with respect to the parameter ϕ , encompassing all model parameters, using the Adam optimizer (Kingma & Ba, 2014). Notably, the additional terms introduced do not require any extra user-specified hyperparameters, which is easy to optimize in practice.

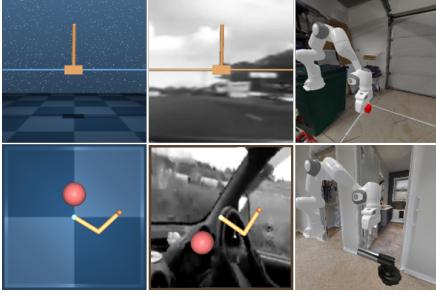


Figure 3. Pixel observations of the DeepMind Control suite (left column) for *cartpole* (top) and *reacher* (bottom), Distracted DeepMind Control suite (middle column) for *cartpole* (top) and *reacher* (bottom), and Mani-skill2 environments with distractions (right column) for *cube* (top) and *faucet* (bottom).

Learning to control With the latent representation and dynamics model, we perform actor-critic policy learning by rolling out trajectories in the latent space. The critic $v_\psi(s_t)$ is trained to predict the discounted cumulative reward given a latent state, and the actor $\pi_\psi(s_t)$ is trained to take the action that maximizes the critic’s prediction, which follows actor-critic training in DreamerV3 (Hafner et al., 2023).

5. Analysis

Our primary goal is to learn good state representations by focusing on two key objectives: latent reconstruction via a masking strategy for compact representations, and employing behavioral similarity for efficient representations. This section will highlight both components are essential for our world model, underscoring their necessity.

Consider an MDP \mathcal{M} as defined in Section 3, with vectorized state variables $\zeta = [s; \xi]$, where $\xi = [\xi^0 \xi^1 \dots \xi^{n-1}]$ is a n -dim vector. We begin with an ideal assumption that our masking strategy only applies on exogenous noise ξ , i.e., $\tilde{\xi} \subseteq \xi$ be an arbitrary subset (a mask) and $\bar{\xi} = \xi \setminus \tilde{\xi}$ be the variables not included in the mask. Then the state reduces to $\bar{\zeta} = [s; \bar{\xi}]$. And we would like to know if the policy $\bar{\pi}$ under reduced MDP $\bar{\mathcal{M}}$ still being optimal for original MDP \mathcal{M} .

Theorem 5.1. If (1) $r(s_t, \xi^i, a_t) = 0 \forall \xi^i \in \bar{\xi}$, (2) $P(s_{t+1}|s_t, \xi, a_t) = P(s_{t+1}|s_t, \bar{\xi}, a_t)$, and (3) $P(\tilde{\xi}_{t+1}, \bar{\xi}_{t+1}|\tilde{\xi}_t, \bar{\xi}_t) = P(\tilde{\xi}_{t+1}|\tilde{\xi}_t) \cdot P(\bar{\xi}_{t+1}|\bar{\xi}_t)$, then we have $\bar{V}_{\bar{\pi}}(\zeta) = V_{\bar{\pi}}(\zeta) \forall \zeta \in \mathcal{Z}$, where $\bar{V}_{\bar{\pi}}(\zeta)$ is the value function under reduced MDP. If $\bar{\pi}$ is optimal for $\bar{\mathcal{M}}$, then $\bar{V}_{\bar{\pi}}(\zeta) = V^*(\zeta) \forall \zeta \in \mathcal{Z}$.

Proof. See Appendix B. \square

It reveals that if we can identify and eliminate exogenous noise without altering the reward or the internal dynamics of the underlying MDP, the resulting value function of this underlying MDP remains optimal with respect to the original problem. This scenario presents an opportunity

for implementing a masking strategy. In practical settings, however, our masking approach involves random patch removal. This randomness does not guarantee the exclusive elimination of exogenous noise. Since elements of the environment crucial to the task may inadvertently be masked, the reward and dynamics can be incorrectly reconstructed, hence the underlying MDP (in latent space) is possibly changed. Consequently, if the masking technique is not sensitive to both the reward and the internal dynamics of the system, an optimal policy can not be assured. This limitation underscores why relying solely on masking-based latent reconstruction is insufficient for learning an effective world model in environments with distractions. Fortunately, the bisimulation principle offers a promising solution. By leveraging this principle, as detailed in Appendix B.2, we can train representations that encapsulate both reward and dynamic information. With bisimulation, the agent can be aware of the reward and the internal dynamics, and therefore can further update towards the optimal policies.

On the other hand, learning state representation only with bisimulation objective is also not sufficient enough for model-based control. In model-based framework, integrating bisimulation-based objective requires to sample consecutive state pairs from an approximate dynamics model, e.g., RSSM in this paper. Though bisimulation objective has practically shown effectiveness in model-free settings (Zhang et al., 2020a; Zang et al., 2022a), (Kemertas & Aumentado-Armstrong, 2021) illustrates that when refer to an approximate dynamics model, this dynamics model needs to meet certain condition to ensure the convergence of the bisimulation principle:

Theorem 5.2. (Kemertas & Aumentado-Armstrong, 2021) Assume \mathcal{S} is compact. For d^π , if the support of an approximate dynamics model \hat{P} , $\text{supp}(\hat{P})$ is a closed subset of \mathcal{S} , then there exists a unique fixed-point d^π , and this metric is bounded: $\text{supp}(\hat{P}) \subseteq \mathcal{S} \Rightarrow \text{diam}(\mathcal{S}; d^\pi) \leq \frac{1}{1-\gamma} (R_{\max} - R_{\min})$.

In practice, the support of an approximate dynamics model cannot be assured to be a subset of the observation space due to the presence of unpredictable exogenous noise. Consequently, when exogenous noise is involved, objectives dependent on transition dynamics, including bisimulation objectives, are likely incapable of filtering out all task-irrelevant information. A numerical counterexample illustrating this point is provided in Appendix B. Therefore, to effectively reduce spatio-temporal redundancy in the observation space, additional methods are necessary. This is the rationale behind our adoption of a masking strategy and latent reconstruction in our approach.

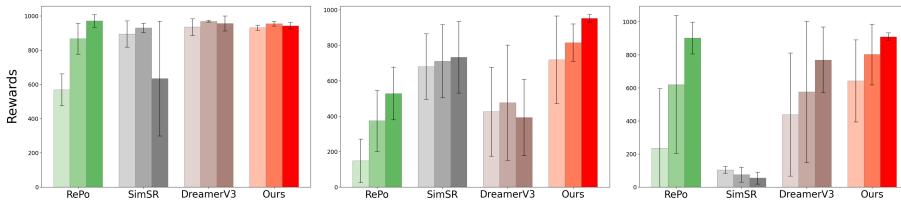


Figure 4. Performance comparison on DMC tasks over 6 seeds in the default setting. Colors from light to dark represent the results evaluated at 100k, 250k, and 500k training steps, respectively, with different colors indicating different models.

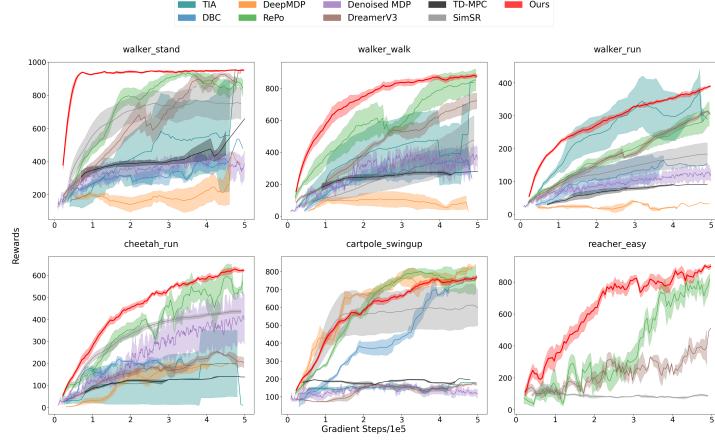


Figure 5. Performance comparison on DMC tasks with one standard error shaded in the distraction setting. The horizontal axis indicates the number of gradient steps. The vertical axis represents the mean return. For our model, RePo, along with DreamerV3 and SimSR, the returns are averaged with six random seeds. For the remaining models, the returns are averaged with three seeds.

6. Experiments

We aim to address the following questions through our experiments : (1) Compared to prior approaches, does our decoder-free model weaken the resulting performance of the policy on downstream tasks? (2) Can we learn effective world models for planning, in the presence of environments containing exogenous spatio-temporal noise structures? (3) We perform ablation studies showing the effectiveness of each of the components in our proposed model (4) Can the proposed Hybrid-RSSM architecture, along with the masking strategy, outperform state-of-the-art Dreamer based models, in presence of exogenous information in data?

Experimental Setup We evaluate our visual image-based continuous control tasks to assess their sample efficiency and overall performance. We perform our experiments in three distinct settings: i) a set of MuJoCo tasks (Todorov et al., 2012) provided by Deepmind Control(DMC) suite (Tassa et al., 2018), ii) a variant of DeepMind Control Suite where the background is replaced with grayscale natural videos from Kinetics dataset (Kay et al., 2017), termed as Distracted DeepMind Control Suite (Zhang et al., 2018), and iii) a benchmark based on the Man-

iskill2 (Gu et al., 2023), enhanced with realistic images of human homes (Chang et al., 2017) as backgrounds and was introduced in (Zhu et al., 2023). Six tasks were tested in the first two settings and two in the last, with a total of 14 tasks. Task examples are depicted in Figure 3.

Baselines We compare our proposed model against leading sample-efficient, model-free and model-based reinforcement learning (RL) methods in continuous control tasks. For model-free mehtods, our baselines include: DBC (Zhang et al., 2020a) and SimSR (Zang et al., 2022a), both of which are two representative bisimulation-based methods. For model-based RL, experimental comparisons are made with TD-MPC (Hansen et al., 2022b), DreamerV3 (Hafner et al., 2023), and its extensions (TIA (Fu et al., 2021), Denoised MDP (Wang et al., 2022a), RePO (Zhu et al., 2023)) that enhance robust representation learning. In this experiment, we use DreamerV3 as our backbone and build on top of it to develop our hybrid structure, where we use an unofficial open-sourced pytorch version of DreamerV3(NM512, 2023). Notably, despite incorporating dual RSSMs, mask branch and raw branch in our framework, our model maintains a slightly smaller overall size compared to the original DreamerV3, which is notable considering the substantial

385
 386 *Table 2.* Summary of performance metrics evaluated at 100K training steps. The performance is quantified in terms of the average score \pm
 387 standard deviation. The highest result for each task is highlighted in bold. For RePo, DreamerV3, and our model, the returns are averaged
 388 with six random seeds. For TIA and Denoised MDP, the returns are averaged with three seeds.
 389
 390

| Task | DreamerV3 | TIA | Denoised MDP | RePo | Ours |
|-------------|---------------|----------------|--------------|-------|---------------|
| Lift Cube | 167±45 | 274±173 | 155±71 | 83±22 | 284±85 |
| Turn Faucet | 138±83 | 47±23 | 71±27 | 92±88 | 278±97 |

392 size of the decoder parameters in DreamerV3. Detailed
 393 descriptions of our model are provided in the Appendix D.1.
 394

395 **Results on DMC tasks with default settings.** As shown
 396 in the left column of Figure 3, the default setting, which is
 397 provided by DMC, has simple backgrounds for the pixel
 398 observations. Figure 4 shows that our model consistently
 400 surpasses all baselines including RePo at 100k, 250k and
 401 500k training steps in all three tasks, showcasing superior
 402 sample efficiency and final performance. Our model
 403 consistently equals or betters the performance of DreamerV3,
 404 illustrating our robustness against performance loss from
 405 omitting pixel-level reconstruction. This also highlights that
 406 the performance improvements of our model are primarily
 407 attributed to the innovative hybrid RSSM structure and
 408 objectives, rather than an increase in size.

409 **Results on DMC tasks with distraction settings.** Figure 5
 410 illustrates our model’s ability to ignore irrelevant information,
 411 outperforming most other models in various tasks. This
 412 underscores our method’s resilience and efficiency in learning
 413 exo-free policies even in the presence of significant
 414 distractor information. Notably, we have almost the lowest
 415 variance across all tasks, which illustrates the robustness
 416 of our hybrid architecture, showing that our HRSSM is
 417 well-suited for model-based agents and is capable of learning
 418 compact and effective representations and dynamics.
 419 However, in the *cartpole_swingup* task, our model slightly
 420 underperforms compared to DeepMDP and RePo. This may
 421 be due to our random masking strategy, which might inad-
 422 vertently hide crucial elements like the small pole, crucial
 423 for task-relevant information. A learned masking strategy
 424 could be more effective than random masking in such cases,
 425 which is deserved to further investigation.

426 **Realistic Maniskill** Table 2 not only demonstrates the com-
 427 petitive performance of our method but also underlines its
 428 distinct advantages in terms of consistency and robustness
 429 across different tasks. In the *Lift Cube* task, our method
 430 achieved a competitive score of 274 ± 35 , paralleling TIA.
 431 However, the significantly lower variance in our results indi-
 432 cates superior consistency and reliability. This is critical in
 433 real-world scenarios where predictability and stability are
 434 as crucial as performance. In *Turn Faucet*, our method’s
 435 superiority is even more pronounced, substantially higher
 436 than its closest competitor. This not only showcases our
 437 method’s ability to handle complex tasks efficiently but also
 438 its robust state representation.

392 **Ablation Studies** Our model comprises two key elements:
 393 mask-based latent reconstruction and a similarity objective
 394 guided by the bisimulation principle. We present their em-
 395 pirical impacts in the distraction setting of DMC tasks in
 396 Appendix E.3. To evaluate mask-based latent reconstruc-
 397 tion, we eliminated the mask branch and reverted our hy-
 398 brid RSSM to a standard RSSM, also omitting the cube
 399 masking and the latent reconstruction loss. For the bisim-
 400 ulation principle ablation, we simply removed the similarity
 401 loss. Results indicate that models lacking these components
 402 underperform relative to the full model, showcasing their
 403 critical importance in our framework.

404 **Training Time Comparison** As our hybrid structure incor-
 405 porates two RSSMs, one might wonder about the computa-
 406 tional efficiency of our framework. Notably, gradients are
 407 only backpropagated through the mask branch, while the
 408 parameters of the RSSM in the raw branch are updated via
 409 Exponential Moving Average (EMA). Moreover, since we
 410 utilize the same historical representation, the computational
 411 time required for the forward process is considerably less
 412 than twice as much. To validate this, we compared the wall-
 413 clock training time of our method against DreamerV3, with
 414 the results provided in Appendix E.2. These results confirm
 415 that our method is comparable to the original DreamerV3
 416 in terms of computational efficiency, without incurring sub-
 417 stantial additional time costs.

7. Discussion

418 **Limitations and Future Work** Our approach’s potential
 419 limitation lies in the lack of a task-specific masking strategy,
 420 which could partially damage the endogenous state and
 421 slightly reduce the final performance. Future improvements
 422 could involve signal-to-noise ratios (Tomar et al., 2023) to
 423 reduce the original image, aiming to identify the minimal
 424 information essential for the task.

425 **Conclusion:** In this paper, we presented a new framework
 426 to learn state representations and dynamics in the presence
 427 of exogenous noise. We introduced the masking strategy and
 428 latent reconstruction to eliminate redundant spatio-temporal
 429 information, and employed bisimulation principle to capture
 430 task-relevant information. Addressing co-training instabili-
 431 ties, we further developed a hybrid RSSM structure. Empirical
 432 results demonstrated the effectiveness of our model.

440 8. Broader Impact

441 This paper synthesizes theoretical and empirical results to
 442 build more capable Model-Based Reinforcement Learning
 443 (MBRL) agents in settings with exogenous noise. In real-
 444 world applications, distractions are prevalent across different
 445 scenarios. Enabling model-based agents to learn control
 446 from such scenarios can be beneficial not only for solving
 447 complex tasks but also for increasing sample efficiency during
 448 deployments in environments with varying contexts. Our
 449 framework is not only theoretically sound but also technically
 450 straightforward to implement and empirically competitive.
 451 We believe that developing MBRL agents by focusing
 452 on the compactness and effectiveness of the representation
 453 and dynamics is an important step towards creating more
 454 applicable Artificial General Intelligence (AGI).

456 References

457 Agarwal, A., Kakade, S., Krishnamurthy, A., and Sun, W.
 458 Flambe: Structural complexity and representation learning
 459 of low rank mdps. *Advances in neural information*
 460 *processing systems*, 33:20095–20107, 2020.

461 Agarwal, R., Schwarzer, M., Castro, P. S., Courville, A. C.,
 462 and Bellemare, M. G. Deep reinforcement learning at
 463 the edge of the statistical precipice. In Ranzato, M.,
 464 Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan,
 465 J. W. (eds.), *Advances in Neural Information Processing*
 466 *Systems 34: Annual Conference on Neural Information*
 467 *Processing Systems 2021, NeurIPS 2021, December 6-14,*
 468 *2021, virtual*, pp. 29304–29320, 2021.

469 Asadi, K., Misra, D., Kim, S., and Littman, M. L. Combating
 470 the compounding-error problem with a multi-step
 471 model. *arXiv preprint arXiv:1905.13320*, 2019.

472 Burda, Y., Edwards, H., Storkey, A., and Klimov, O. Exploration
 473 by random network distillation. *arXiv preprint*
 474 *arXiv:1810.12894*, 2018.

475 Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P.,
 476 and Joulin, A. Unsupervised learning of visual features
 477 by contrasting cluster assignments. *Advances in neural*
 478 *information processing systems*, 33:9912–9924, 2020.

479 Castro, P. S. Scalable methods for computing state similarity
 480 in deterministic markov decision processes. In *Proceedings*
 481 *of the AAAI Conference on Artificial Intelligence*,
 482 volume 34, pp. 10069–10076, 2020.

483 Castro, P. S., Kastner, T., Panangaden, P., and Rowland,
 484 M. Mico: Improved representations via sampling-based
 485 state similarity for markov decision processes. *Advances*
 486 *in Neural Information Processing Systems*, 34:30113–
 487 30126, 2021.

488 Chang, A., Dai, A., Funkhouser, T., Halber, M., Niessner,
 489 M., Savva, M., Song, S., Zeng, A., and Zhang, Y. Matterport3d:
 490 Learning from rgb-d data in indoor environments.
arXiv preprint arXiv:1709.06158, 2017.

491 Chitnis, R. and Lozano-Pérez, T. Learning compact models
 492 for planning with exogenous processes. In *Conference*
 493 *on Robot Learning*, pp. 813–822. PMLR, 2020.

494 Cuturi, M. Sinkhorn distances: Lightspeed computation
 495 of optimal transport. *Advances in neural information*
 496 *processing systems*, 26, 2013.

497 Deng, F., Jang, I., and Ahn, S. Dreamerpro: Reconstruction-
 498 free model-based reinforcement learning with prototypical
 499 representations. In *International Conference on Machine*
 500 *Learning*, pp. 4956–4975. PMLR, 2022.

501 Efroni, Y., Misra, D., Krishnamurthy, A., Agarwal, A., and
 502 Langford, J. Provably filtering exogenous distractors
 503 using multistep inverse dynamics. In *International Conference*
 504 *on Learning Representations*, 2021.

505 Efroni, Y., Foster, D. J., Misra, D., Krishnamurthy, A., and
 506 Langford, J. Sample-efficient reinforcement learning in
 507 the presence of exogenous information. In *Conference*
 508 *on Learning Theory*, pp. 5062–5127. PMLR, 2022.

509 Feichtenhofer, C., Fan, H., Li, Y., and He, K. Masked
 510 autoencoders as spatiotemporal learners. *CoRR*,
 511 abs/2205.09113, 2022. doi: 10.48550/arXiv.2205.
 512 09113. URL <https://doi.org/10.48550/arXiv.2205.09113>.

513 Ferns, N., Panangaden, P., and Precup, D. Bisimulation
 514 metrics for continuous markov decision processes. *SIAM*
 515 *Journal on Computing*, 40(6):1662–1714, 2011.

516 Ferns, N., Castro, P. S., Precup, D., and Panangaden, P.
 517 Methods for computing state similarity in markov deci-
 518 sion processes. *arXiv preprint arXiv:1206.6836*, 2012a.

519 Ferns, N., Panangaden, P., and Precup, D. Metrics
 520 for finite markov decision processes. *arXiv preprint*
arXiv:1207.4114, 2012b.

521 Fu, X., Yang, G., Agrawal, P., and Jaakkola, T. Learning
 522 task informed abstractions. In *International Conference*
 523 *on Machine Learning*, pp. 3480–3491. PMLR, 2021.

524 Gelada, C., Kumar, S., Buckman, J., Nachum, O., and Bel-
 525 lemare, M. G. Deepmdp: Learning continuous latent space
 526 models for representation learning. In *International Con-
 527 ference on Machine Learning*, pp. 2170–2179. PMLR,
 528 2019.

529 Givan, R., Dean, T. L., and Greig, M. Equivalence notions
 530 and model minimization in markov decision processes.
Artif. Intell., 147(1-2):163–223, 2003.

- 495 Gu, J., Xiang, F., Li, X., Ling, Z., Liu, X., Mu, T., Tang,
 496 Y., Tao, S., Wei, X., Yao, Y., et al. Maniskill2: A unified
 497 benchmark for generalizable manipulation skills. *arXiv*
 498 preprint arXiv:2302.04659, 2023.
- 500 Ha, D. and Schmidhuber, J. Recurrent world models facil-
 501 itate policy evolution. *Advances in neural information*
 502 *processing systems*, 31, 2018.
- 503 Hafner, D., Lillicrap, T., Ba, J., and Norouzi, M. Dream to
 504 control: Learning behaviors by latent imagination. *arXiv*
 505 preprint arXiv:1912.01603, 2019a.
- 506 Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D.,
 507 Lee, H., and Davidson, J. Learning latent dynamics for
 508 planning from pixels. In *International conference on*
 509 *machine learning*, pp. 2555–2565. PMLR, 2019b.
- 510 Hafner, D., Lillicrap, T., Norouzi, M., and Ba, J. Mastering
 511 atari with discrete world models. *arXiv preprint*
 512 arXiv:2010.02193, 2020.
- 513 Hafner, D., Pasukonis, J., Ba, J., and Lillicrap, T. Mastering
 514 diverse domains through world models. *arXiv preprint*
 515 arXiv:2301.04104, 2023.
- 516 Hamrick, J. B. Analogues of mental simulation and imagi-
 517 nation in deep learning. *Current Opinion in Behavioral*
 518 *Sciences*, 29:8–16, 2019.
- 519 Hansen, N. and Wang, X. Generalization in reinforcement
 520 learning by soft data augmentation. In *2021 IEEE Inter-*
 521 *national Conference on Robotics and Automation (ICRA)*,
 522 pp. 13611–13617. IEEE, 2021.
- 523 Hansen, N., Su, H., and Wang, X. Stabilizing deep q-
 524 learning with convnets and vision transformers under
 525 data augmentation. *Advances in neural information pro-*
 526 *cessing systems*, 34:3680–3693, 2021.
- 527 Hansen, N., Lin, Y., Su, H., Wang, X., Kumar, V., and
 528 Rajeswaran, A. Modem: Accelerating visual model-
 529 based reinforcement learning with demonstrations. *arXiv*
 530 preprint arXiv:2212.05698, 2022a.
- 531 Hansen, N., Wang, X., and Su, H. Temporal difference
 532 learning for model predictive control. *arXiv preprint*
 533 arXiv:2203.04955, 2022b.
- 534 Hansen, N., Su, H., and Wang, X. Td-mpc2: Scalable, ro-
 535 bust world models for continuous control. *arXiv preprint*
 536 arXiv:2310.16828, 2023.
- 537 He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick,
 538 R. Masked autoencoders are scalable vision learners. In
 539 *Proceedings of the IEEE/CVF conference on computer*
 540 *vision and pattern recognition*, pp. 16000–16009, 2022.
- 541 Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup,
 542 D., and Meger, D. Deep reinforcement learning that mat-
 543 ters. In McIlraith, S. A. and Weinberger, K. Q. (eds.),
 544 *Proceedings of the Thirty-Second AAAI Conference on*
 545 *Artificial Intelligence, (AAAI-18), the 30th innovative Ap-*
 546 *plications of Artificial Intelligence (IAAI-18), and the*
 547 *8th AAAI Symposium on Educational Advances in Arti-*
 548 *ficial Intelligence (EAAI-18), New Orleans, Louisiana,*
 549 *USA, February 2-7, 2018*, pp. 3207–3214. AAAI Press,
 550 2018. doi: 10.1609/AAAI.V32I1.11694. URL <https://doi.org/10.1609/aaai.v32i1.11694>.
- 551 Islam, R., Tomar, M., Lamb, A., Efroni, Y., Zang, H., Di-
 552 dolkar, A., Misra, D., Li, X., van Seijen, H., Combes, R.
 553 T. d., et al. Agent-controller representations: Principled
 554 offline rl with rich exogenous information. *arXiv preprint*
 555 arXiv:2211.00164, 2022.
- 556 Islam, R., Tomar, M., Lamb, A., Efroni, Y., Zang, H., Didolkar,
 557 A. R., Misra, D., Li, X., Van Seijen, H., Des Combes, R. T.,
 558 et al. Principled offline rl in the presence of rich exogenous
 559 information. 2023.
- 560 Janner, M., Fu, J., Zhang, M., and Levine, S. When to trust
 561 your model: Model-based policy optimization. *Advances*
 562 *in neural information processing systems*, 32, 2019.
- 563 Janner, M., Mordatch, I., and Levine, S. gamma-models:
 564 Generative temporal difference learning for infinite-
 565 horizon prediction. *Advances in Neural Information Pro-*
 566 *cessing Systems*, 33:1724–1735, 2020.
- 567 Kalweit, G. and Boedecker, J. Uncertainty-driven imagi-
 568 nation for continuous deep reinforcement learning. In
 569 *Conference on Robot Learning*, pp. 195–206. PMLR,
 570 2017.
- 571 Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier,
 572 C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T.,
 573 Natsev, P., et al. The kinetics human action video dataset.
 574 *arXiv preprint arXiv:1705.06950*, 2017.
- 575 Kemertas, M. and Aumentado-Armstrong, T. Towards ro-
 576 bust bisimulation metric learning. *Advances in Neural*
 577 *Information Processing Systems*, 34:4764–4777, 2021.
- 578 Kingma, D. P. and Ba, J. Adam: A method for stochastic
 579 optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- 580 Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X.,
 581 Sutskever, I., and Welling, M. Improved variational infer-
 582 ence with inverse autoregressive flow. *Advances in neural*
 583 *information processing systems*, 29, 2016.
- 584 Kostrikov, I., Yarats, D., and Fergus, R. Image augmentation
 585 is all you need: Regularizing deep reinforcement learning
 586 from pixels. *arXiv preprint arXiv:2004.13649*, 2020.

- 550 Lamb, A., Islam, R., Efroni, Y., Didolkar, A. R., Misra,
 551 D., Foster, D. J., Molu, L. P., Chari, R., Krishnamurthy,
 552 A., and Langford, J. Guaranteed discovery of control-
 553 endogenous latent states with multi-step inverse models.
 554 *Transactions on Machine Learning Research*, 2022.
- 555 Larsen, K. G. and Skou, A. Bisimulation through proba-
 556 bility testing. In *Conference Record of the Sixteenth*
 557 *Annual ACM Symposium on Principles of Programming*
 558 *Languages, Austin, Texas, USA, January 11-13, 1989*, pp.
 559 344–352. ACM Press, 1989.
- 560 Laskin, M., Srinivas, A., and Abbeel, P. Curl: Contrastive
 561 unsupervised representations for reinforcement learning.
 562 In *International Conference on Machine Learning*, pp.
 563 5639–5650. PMLR, 2020.
- 564 Liu, F., Liu, H., Grover, A., and Abbeel, P. Masked autoen-
 565 coding for scalable and generalizable decision making.
 566 *Advances in Neural Information Processing Systems*, 35:
 567 12608–12618, 2022.
- 568 NM512. Dreamerv3 pytorch implementation. <https://github.com/NM512/dreamerv3-torch>, 2023.
- 569 Okada, M. and Taniguchi, T. Dreaming: Model-based re-
 570 inforcement learning by latent imagination without re-
 571 construction. In *2021 ieee international conference on*
 572 *robotics and automation (icra)*, pp. 4209–4215. IEEE,
 573 2021.
- 574 Poudel, R. P., Pandya, H., Liwicki, S., and Cipolla, R.
 575 Recore: Regularized contrastive representation learning
 576 of world model. *arXiv preprint arXiv:2312.09056*, 2023.
- 577 Rafailov, R., Yu, T., Rajeswaran, A., and Finn, C. Offline
 578 reinforcement learning from images with latent space
 579 models. In *Learning for Dynamics and Control*, pp. 1154–
 580 1168. PMLR, 2021.
- 581 Ramakrishnan, S. K., Gokaslan, A., Wijmans, E.,
 582 Maksymets, O., Clegg, A., Turner, J., Undersander,
 583 E., Galuba, W., Westbury, A., Chang, A. X., et al.
 584 Habitat-matterport 3d dataset (hm3d): 1000 large-scale
 585 3d environments for embodied ai. *arXiv preprint*
 586 *arXiv:2109.08238*, 2021.
- 587 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R.,
 588 Parikh, D., and Batra, D. Grad-cam: Visual explana-
 589 tions from deep networks via gradient-based localiza-
 590 tion. In *Proceedings of the IEEE international conference on*
 591 *computer vision*, pp. 618–626, 2017.
- 592 Seo, Y., Hafner, D., Liu, H., Liu, F., James, S., Lee, K., and
 593 Abbeel, P. Masked world models for visual control. In
 594 *Conference on Robot Learning*, pp. 1332–1344. PMLR,
 595 2023a.
- 596 Seo, Y., Kim, J., James, S., Lee, K., Shin, J., and Abbeel,
 597 P. Multi-view masked world models for visual robotic
 598 manipulation. *arXiv preprint arXiv:2302.02408*, 2023b.
- 599 Stone, A., Ramirez, O., Konolige, K., and Jonschkowski,
 600 R. The distracting control suite - A challenging bench-
 601 mark for reinforcement learning from pixels. *CoRR*,
 602 abs/2101.02722, 2021. URL <https://arxiv.org/abs/2101.02722>.
- 603 Sutton, R. S. Integrated architectures for learning, plan-
 604 ning, and reacting based on approximating dynamic pro-
 605 gramming. In *Machine learning proceedings 1990*, pp.
 606 216–224. Elsevier, 1990.
- 607 Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.
 608 d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq,
 609 A., et al. Deepmind control suite. *arXiv preprint*
 610 *arXiv:1801.00690*, 2018.
- 611 Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics
 612 engine for model-based control. In *2012 IEEE/RSJ interna-*
613 tional conference on intelligent robots and systems, pp.
 614 5026–5033. IEEE, 2012.
- 615 Tomar, M., Islam, R., Levine, S., and Bachman, P. Ignorance
 616 is bliss: Robust control via information gating. *arXiv*
 617 *preprint arXiv:2303.06121*, 2023.
- 618 Tong, Z., Song, Y., Wang, J., and Wang, L. Videomae:
 619 Masked autoencoders are data-efficient learners for self-
 620 supervised video pre-training. *Advances in neural infor-*
621 mation processing systems, 35:10078–10093, 2022.
- 622 Wang, T., Du, S. S., Torralba, A., Isola, P., Zhang, A., and
 623 Tian, Y. Denoised mdps: Learning world models better
 624 than the world itself. *arXiv preprint arXiv:2206.15477*,
 625 2022a.
- 626 Wang, Z., Xiao, X., Xu, Z., Zhu, Y., and Stone, P. Causal
 627 dynamics learning for task-independent state abstraction.
 628 *arXiv preprint arXiv:2206.13452*, 2022b.
- 629 Wei, C., Fan, H., Xie, S., Wu, C.-Y., Yuille, A., and Feicht-
 630 enhofer, C. Masked feature prediction for self-supervised
 631 visual pre-training. In *Proceedings of the IEEE/CVF Con-*
632 ference on Computer Vision and Pattern Recognition, pp.
 633 14668–14678, 2022.
- 634 Williams, G., Aldrich, A., and Theodorou, E. A. Model
 635 predictive path integral control: From theory to parallel
 636 computation. *Journal of Guidance, Control, and Dynam-
 637 ics*, 40(2):344–357, 2017.
- 638 Xiao, C., Wu, Y., Ma, C., Schuurmans, D., and Müller, M.
 639 Learning to combat compounding-error in model-based
 640 reinforcement learning. *arXiv preprint arXiv:1912.11206*,
 641 2019.

- 605 Yarats, D., Fergus, R., Lazaric, A., and Pinto, L. Mastering
606 visual continuous control: Improved data-augmented re-
607inforcement learning. *arXiv preprint arXiv:2107.09645*,
608 2021a.
- 609 Yarats, D., Zhang, A., Kostrikov, I., Amos, B., Pineau, J.,
610 and Fergus, R. Improving sample efficiency in model-
611 free reinforcement learning from images. In *Proceedings
612 of the AAAI Conference on Artificial Intelligence*, vol-
613 ume 35, pp. 10674–10681, 2021b.
- 614
- 615 Yu, T., Zhang, Z., Lan, C., Lu, Y., and Chen, Z. Mask-
616 based latent reconstruction for reinforcement learning.
617 *Advances in Neural Information Processing Systems*, 35:
618 25117–25131, 2022.
- 619
- 620 Zang, H., Li, X., and Wang, M. Simsr: Simple distance-
621 based state representations for deep reinforcement learn-
622 ing. In *Proceedings of the AAAI Conference on Artificial
623 Intelligence*, volume 36, pp. 8997–9005, 2022a.
- 624
- 625 Zang, H., Li, X., Yu, J., Liu, C., Islam, R., Combes, R.
626 T. D., and Laroche, R. Behavior prior representation
627 learning for offline reinforcement learning. *arXiv preprint
628 arXiv:2211.00863*, 2022b.
- 629
- 630 Zhang, A., Wu, Y., and Pineau, J. Natural environment
631 benchmarks for reinforcement learning. *arXiv preprint
632 arXiv:1811.06032*, 2018.
- 633
- 634 Zhang, A., Lyle, C., Sodhani, S., Filos, A., Kwiatkowska,
635 M., Pineau, J., Gal, Y., and Precup, D. Invariant causal
636 prediction for block mdps. In *International Conference
637 on Machine Learning*, pp. 11214–11224. PMLR, 2020a.
- 638
- 639 Zhang, A., McAllister, R., Calandra, R., Gal, Y., and
640 Levine, S. Learning invariant representations for rein-
641 forcement learning without reconstruction. *arXiv preprint
642 arXiv:2006.10742*, 2020b.
- 643
- 644 Zhao, Y., Zhao, W., Boney, R., Kannala, J., and Pajarin, J.
645 Simplified temporal consistency reinforcement learning.
646 *arXiv preprint arXiv:2306.09466*, 2023.
- 647
- 648 Zhu, C., Simchowitz, M., Gadipudi, S., and Gupta, A. Repo:
649 Resilient model-based reinforcement learning by regu-
650 larizing posterior predictability. *CoRR*, abs/2309.00082,
651 2023. doi: 10.48550/ARXIV.2309.00082. URL <https://doi.org/10.48550/arXiv.2309.00082>.
- 652
- 653
- 654
- 655
- 656
- 657
- 658
- 659

A. Hyperparameters

We present all hyperparameters in Table 3.

| Name | Symbol | Value |
|----------------------------|------------------|--|
| General | | |
| Replay capacity (FIFO) | — | 10^6 |
| Batch size | B | 16 |
| Batch length | T | 64 |
| Activation | — | LayerNorm + SiLU |
| World Model | | |
| Number of latents | — | 32 |
| Classes per latent | — | 32 |
| Learning rate | — | 10^{-4} |
| Adam epsilon | ϵ | 10^{-8} |
| Gradient clipping | — | 1000 |
| Actor Critic | | |
| Imagination horizon | H | 15 |
| Discount horizon | $1/(1 - \gamma)$ | 333 |
| Return lambda | λ | 0.95 |
| Critic EMA decay | — | 0.98 |
| Critic EMA regularizer | — | 1 |
| Return normalization scale | S | $\text{Per}(R, 95) - \text{Per}(R, 5)$ |
| Return normalization limit | L | 1 |
| Return normalization decay | — | 0.99 |
| Actor entropy scale | η | $3 \cdot 10^{-4}$ |
| Learning rate | — | $3 \cdot 10^{-5}$ |
| Adam epsilon | ϵ | 10^{-5} |
| Gradient clipping | — | 100 |
| Masking | | |
| Mask ratio | — | 50% |
| Cube spatial size | $h \times w$ | 10×10 |
| Cube depth | k | 4 |

Table 3. Our model’s hyperparameters, which are the same across all tasks in DMControl and Realistic Maniskill.

B. Analysis and Example

B.1. Masking Strategy

Consider an MDP \mathcal{M} as defined in Section 3, with vectorized state variables $\zeta = [s; \xi]$, where $\xi = [\xi^0 \xi^1 \dots \xi^n]$ is a n -dim vector. We begin with an ideal assumption that our masking strategy only applies on exogenous noise ξ , i.e., $\tilde{\xi} \subseteq \xi$ be an arbitrary subset (a mask) and $\bar{\xi} = \xi \setminus \tilde{\xi}$ be the variables not included in the mask. Then the state reduces to $\bar{\zeta} = [s; \bar{\xi}]$. And we would like to know if the policy $\bar{\pi}$ under reduced MDP $\bar{\mathcal{M}}$ still being optimal for original MDP \mathcal{M} .

Theorem B.1. If (1) $r(s_t, \xi_t^i, a_t) = 0 \forall \xi^i \in \bar{\xi}$, (2) $P(s_{t+1}|s_t, \xi, a_t) = P(s_{t+1}|s_t, \bar{\xi}, a_t)$, and (3) $P(\tilde{\xi}_{t+1}, \bar{\xi}_{t+1}|\tilde{\xi}_t, \bar{\xi}_t) = P(\tilde{\xi}_{t+1}|\tilde{\xi}_t) \cdot P(\bar{\xi}_{t+1}|\bar{\xi}_t)$, then we have $\bar{V}_{\bar{\pi}}(\bar{\zeta}) = V_{\bar{\pi}}(\zeta) \forall \zeta \in \mathcal{Z}$, where $\bar{V}_{\bar{\pi}}(\bar{\zeta})$ is the value function under reduced MDP. If $\bar{\pi}$ is optimal for $\bar{\mathcal{M}}$, then $\bar{V}_{\bar{\pi}}(\bar{\zeta}) = V^*(\zeta) \forall \zeta \in \mathcal{Z}$.

Proof. This proof mimics the proof of Theorem 1 in (Chitnis & Lozano-Pérez, 2020). Consider an arbitrary state $\zeta \in \mathcal{Z}$,

715 and its reduced state $\bar{\zeta}$, we have the following equations:

$$V_{\bar{\pi}}(\zeta) = R(\zeta, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\zeta'} P(\zeta' | \zeta, \bar{\pi}(\bar{\zeta})) \cdot V_{\bar{\pi}}(\zeta'). \quad (8)$$

$$\bar{V}_{\bar{\pi}}(\bar{\zeta}) = R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\bar{\zeta}'} P(\bar{\zeta}' | \bar{\zeta}, \bar{\pi}(\bar{\zeta})) \cdot V_{\bar{\pi}}(\bar{\zeta}'). \quad (9)$$

722 Now suppose $V_{\bar{\pi}}^k(\zeta) = \bar{V}_{\bar{\pi}}^k(\bar{\zeta}) \forall \zeta \in \mathcal{Z}$, for some k .

$$\begin{aligned} V_{\bar{\pi}}^{k+1}(\zeta) &= R(\zeta, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\zeta'} P(\zeta' | \zeta, \bar{\pi}(\bar{\zeta})) \cdot V_{\bar{\pi}}^k(\zeta') \\ &= R(s, \bar{\pi}(\bar{\zeta})) + \sum_{i=1}^n R^i(s, \xi^i, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\zeta'} P(s' | s, \bar{\pi}(\bar{\zeta}), \xi) \cdot P(\xi' | \xi) \cdot V_{\bar{\pi}}^k(\zeta') \\ &= R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\zeta'} P(s' | s, \bar{\pi}(\bar{\zeta}), \xi) \cdot P(\xi' | \xi) \cdot V_{\bar{\pi}}^k(\zeta') \\ &= R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\zeta'} P(s' | s, \bar{\pi}(\bar{\zeta}), \bar{\xi}) \cdot P(\xi' | \xi) \cdot V_{\bar{\pi}}^k(\zeta') \\ &= R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{s', \bar{\xi}', \tilde{\xi}'} P(s' | s, \bar{\pi}(\bar{\zeta}), \bar{\xi}) \cdot P(\tilde{\xi}', \bar{\xi}' | \tilde{\xi}, \bar{\xi}) \cdot V_{\bar{\pi}}^k(\zeta') \\ &= R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{s', \bar{\xi}', \tilde{\xi}'} P(s' | s, \bar{\pi}(\bar{\zeta}), \bar{\xi}) \cdot P(\tilde{\xi}' | \tilde{\xi}) P(\bar{\xi}' | \bar{\xi}) \cdot V_{\bar{\pi}}^k(\zeta') \\ &= R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{s', \bar{\xi}', \tilde{\xi}'} P(s' | s, \bar{\pi}(\bar{\zeta}), \bar{\xi}) \cdot P(\tilde{\xi}' | \tilde{\xi}) P(\bar{\xi}' | \bar{\xi}) \cdot V_{\bar{\pi}}^k(s', \bar{\xi}') \\ &= R(\bar{\zeta}, \bar{\pi}(\bar{\zeta})) + \gamma \sum_{\bar{\zeta}} P(\bar{\zeta}' | \bar{\zeta}, \bar{\pi}(\bar{\zeta})) \cdot V_{\bar{\pi}}^k(\bar{\zeta}'). \\ &= \bar{V}_{\bar{\pi}}^{k+1}(\bar{\zeta}). \end{aligned}$$

746 Therefore, we have that $\bar{V}_{\bar{\pi}}(\bar{\zeta}) = V_{\bar{\pi}}(\zeta) \forall \zeta \in \mathcal{Z}$. And if $\bar{\pi}$ is optimal for $\bar{\mathcal{M}}$, then it is optimal for the full MDP \mathcal{M} as well. \square

B.2. Bisimulation Principle

751 Bisimulation measures equivalence relations on MDPs in a recursive manner: two states are considered equivalent if they
 752 share equivalent distributions over the next equivalent states and have the same immediate reward (Larsen & Skou, 1989;
 753 Givan et al., 2003).

754 **Definition B.2.** Given an MDP \mathcal{M} , an equivalence relation $E \subseteq \mathcal{S} \times \mathcal{S}$ is a bisimulation relation if whenever $(s, u) \in E$
 755 the following properties hold, where \mathcal{S}_E is the state space \mathcal{S} partitioned into equivalence classes defined by E :

- 757 1. $\forall a \in \mathcal{A}, \mathcal{R}(s, a) = \mathcal{R}(u, a)$
- 758 2. $\forall a \in \mathcal{A}, \forall c \in \mathcal{S}_E, \mathcal{P}(s, a)(c) = \mathcal{P}(u, a)(c)$ where $\mathcal{P}(s, a)(c) = \sum_{s' \in c} \mathcal{P}(s, a)(s')$.

761 Two states $s, u \in \mathcal{S}$ are bisimilar if there exists a bisimulation relation E such that $(s, u) \in E$. We denote the largest
 762 bisimulation relation as \sim .

764 However, bisimulation, by considering equivalence for all actions including bad ones, often leads to "pessimistic" outcomes.
 765 To address this, (Castro, 2020) introduced π -bisimulation, which eliminates the need to consider every action and instead
 766 focuses on actions induced by a policy π .

767 **Definition B.3.** (Castro, 2020) Given an MDP \mathcal{M} , an equivalence relation $E^\pi \subseteq \mathcal{S} \times \mathcal{S}$ is a π -bisimulation relation if the
 768 following properties hold whenever $(s, u) \in E^\pi$:

770 1. $r(s, \pi) = r(u, \pi)$

771 2. $\forall C \in \mathcal{S}_{E^\pi}, T(C|s, \pi) = T(C|u, \pi)$

773
774 where \mathcal{S}_{E^π} is the state space \mathcal{S} partitioned into equivalence classes defined by E^π . Two states $s, u \in \mathcal{S}$ are π -bisimilar if
775 there exists a π -bisimulation relation E^π such that $(s, u) \in E^\pi$.

776 However, π -bisimulation is still too strict to be practically applied at scale, as it treats equivalence as a binary property:
777 either two states are equivalent or not, making it highly sensitive to perturbations in numerical values of model parameters.
778 This issue becomes even more pronounced when deep frameworks are employed. To address this, (Castro, 2020) further
780 proposed a π -bisimulation metric that incorporates the absolute difference between immediate rewards of two states and the
781 1-Wasserstein distance (\mathcal{W}_1) between the transition distributions conditioned on the two states and the policy π :

782 **Theorem B.4** ((Castro, 2020)). Define $\mathcal{F}^\pi : \mathcal{M} \rightarrow \mathcal{M}$ by $\mathcal{F}^\pi(d)(u, v) = |R(u, \pi) - R(v, \pi)| + \gamma\mathcal{W}_1(d)(P_u^\pi, P_v^\pi)$, then
783 \mathcal{F}^π has a least fixed point d_\sim^π , and d_\sim^π is a π -bisimulation metric.

784 It suffices to show that above fixed-point updates are contraction mappings. Then the existence of a unique metric can be
785 proved by invoke the Banach fixed-point theorem (Ferns et al., 2011). An essential assumption is that the state space \mathcal{S}
786 should be compact¹. And the compactness of \mathcal{S} implies that the metric space over this state space is complete such that the
787 Banach fixed-point theorem can be applied. And when considering the approximate dynamics, the situation becomes more
788 complicated. (Kemertas & Aumentado-Armstrong, 2021) show that:

789 **Theorem B.5** ((Kemertas & Aumentado-Armstrong, 2021)). Assume \mathcal{S} is compact. For d^π , if the support of an approximate
790 dynamics model \hat{P} , $\text{supp}(\hat{P})$ is a closed subset of \mathcal{S} , then there exists a unique fixed-point d^π , and this metric is bounded:

791
$$\text{supp}(\hat{P}) \subseteq \mathcal{S} \Rightarrow \text{diam}(\mathcal{S}; d^\pi) \leq \frac{1}{1-\gamma} (R_{\max} - R_{\min}) \quad (10)$$

792 *Proof.* The proof adapts from (Kemertas & Aumentado-Armstrong, 2021), which is also a slight generalization of the
793 distance bounds given in Theorem 3.12 of (Ferns et al., 2011).

794
$$d^\pi(u, v) = |R(u, \pi) - R(v, \pi)| + \gamma W(d)(P(\cdot|u, \pi), P(\cdot|v, \pi)) \leq R_{\max} - R_{\min} + \gamma \text{diam}(\mathcal{S}; d^\pi), \forall (u, v) \in \mathcal{S} \times \mathcal{S}, \quad (11)$$

795 with the use of Lemma 5 in (Kemertas & Aumentado-Armstrong, 2021), we have:

796
$$\begin{aligned} \text{diam}(\mathcal{S}; d^\pi) &\leq R_{\max} - R_{\min} + \gamma \text{diam}(\mathcal{S}; d^\pi) \\ &\leq \frac{1}{1-\gamma} (R_{\max} - R_{\min}) \end{aligned} \quad (12)$$

800 \square

801 In this paper, our bisimulation objective is defined as follows:

802
$$\mathcal{F}^\pi d(s_i, s_j) = |r_{s_i}^\pi - r_{s_j}^\pi| + \gamma E_{\substack{s_{i+1} \sim \hat{P}_{s_i}^a, \\ s_{j+1} \sim \hat{P}_{s_j}^a}} [d(s_{i+1}, s_{j+1})], \quad (13)$$

803 where we sample the next state pairs from an approximated dynamics model RSSM instead of the ground-truth dynamics,
804 and use the independent coupling instead of computing Wasserstein distance. In principle, iteration on conventional state
805 space is acceptable with such a method. While in practice, the above requirement is hard to be satisfied as we learn state
806 representation from an noisy observation space that includes unpredictable exogenous noise.

807 B.3. Counterexample

808 Consider two vectorized states $u = (1, 2, 3, 1, 1)$, $v = (2, 1, 1, 1, 1)$, where the last two dimension of these states are
809 exogenous noise that irrelevant to the task. Under policy π , their next states are $u' = (2, 2, 1, 1, 1)$, $v' = (1, 1, 2, 1, 1)$

810 ¹A continuous space is compact if and only if it is totally bounded and complete.

825 respectively. Give $\gamma = 0.92$, $r_u^\pi = 0.03$, $r_v^\pi = 0.02$, and with an error of $\epsilon = 0.01$, we almost reach the optimal bisimulation
826 distance:

$$\begin{aligned} d(u, v) &= 0.7955 \\ (r_u^\pi - r_v^\pi) + d(u', v') &= 0.7945 \\ \Delta &= 0.7955 - 0.7945 = 0.001 < \epsilon. \end{aligned} \quad (14)$$

831 Meanwhile, the endogenous states $\bar{u} = (1, 2, 3)$, $\bar{v} = (2, 1, 1)$, also achieve their optimal bisimulation distance:
832

$$\begin{aligned} d(\bar{u}, \bar{v}) &= 0.7638 \\ (r_u^\pi - r_v^\pi) + d(\bar{u}', \bar{v}') &= 0.7612 \\ \Delta &= 0.7638 - 0.7612 = 0.0026 < \epsilon, \end{aligned} \quad (15)$$

833 while u and v still contain exogenous noise. That is to say, only with bisimulation principle is sufficient to learn task-relevant
834 information, while not enough to learn compact representation.
835

840 C. Additional Related Work Discussion

841 In this section, we provide an additional related work description, and a detailed comparison between our model and other
842 baselines that developed based on Dreamer, including TIA, Denoised MDP, DreamerPro, and RePo.
843

844 **State Representation Learning** Recent advancements in Reinforcement Learning (RL) emphasize learning state representations to understand environment structures, with successful methods like CURL (Laskin et al., 2020) and DrQ (Kostrikov et al., 2020; Yarats et al., 2021a) using data augmentation techniques such as cropping and color jittering, yet their efficacy is closely tied to the specific augmentation employed. Approaches like masking-based approaches (Seo et al., 2023b; Yu et al., 2022; Seo et al., 2023a; Liu et al., 2022) aim to reduce spatiotemporal redundancy but often overlook task-relevant information. Bisimulation-based methods (Zhang et al., 2020b; Zang et al., 2022a) focus on learning reward-aware state representations for value-equivalence and sample efficiency, but they face challenges in achieving compact representation spaces since they sample consecutive states from approximated dynamics. Additionally, a branch of research investigates causality to discover causal relationships between state representation and control (Wang et al., 2022b; Lamb et al., 2022; Islam et al., 2023; Efroni et al., 2021; 2022; Fu et al., 2021; Zang et al., 2022b). Our work primarily follows the methods based on bisimulation and masking, while developing a hybrid RSSM structure tailored for model-based agents.
845

846 **TIA (Fu et al., 2021)** extended Dreamer by creating a cooperative two-player game involving two models: the task model and the distractor model. The distractor model aims to disassociate from the reward as much as possible, while the task model focuses on capturing task-relevant information. Both models contribute to a reconstruction process involving an inferred mask in pixel-space. Although TIA shares similarities with our model, such as the use of masks and a dual-model framework, our hybrid RSSM structure differs in that it does not explicitly model exogenous noise, instead employing a random masking strategy. Moreover, our approach has lower time complexity than TIA, as we utilize a shared historical representation for both branches in the framework, eliminating the need for separate gradient computations. While TIA's learned mask effectively removes noise distractors through pixel-wise composition, it falls short in addressing more general noise types, such as observation disturbances caused by a shaky camera. From this perspective, investigating the potential solution of making masking strategy informed from the control task is still worthwhile for many approaches including TIA and ours.
847

848 **Denoised MDP (Wang et al., 2022a)** classified RL information into four types based on controllability and its relevance to rewards, defining useful information as that which is controllable and reward-related. Their approach tends to overlook factors unrelated to control, even if they might influence the reward function. To address this, they introduced a variational mutual information regularizer to separate control and reward-relevant information from overall observations. While this method successfully distinguishes between task-relevant and irrelevant components, Denoised MDP demonstrated higher variance and moderate performance in distraction settings. This may be attributed to its continued reliance on pixel-level reconstruction, which, by focusing on minute details, could unintentionally diminish policy performance in distraction settings. Conversely, our method, eschewing pixel-level reconstruction, flexibly eliminates spatio-temporal redundancies while preserving semantic content, leading to enhanced performance.
849

880 **DreamerPro (Deng et al., 2022)** proposed a reconstruction-free MBRL agent by combining the prototypical representation
 881 learning with temporal dynamics learning. Borrowing idea from SwAV (Caron et al., 2020), they tried to align the temporal
 882 latent state with the cluster assignment of the observation. However, their cluster assignment requires to apply the Sinkhorn
 883 Knopp algorithm (Cuturi, 2013) to update prototypes. This requires more computational cost and more hyperparameters to
 884 tune. Besides, DreamerPro still cannot learn task-relevant information as its representation is not informed by reward.
 885
 886
 887

888 **RePo (Zhu et al., 2023)** developed its representation in a way of maximizing mutual information (MI) between the
 889 current representation and all future rewards while minimizing the mutual information between the representation and
 890 observation. Excluding pixel-level reconstruction, they ensure latents predictable by optimizing a variational lower bound on
 891 the MI-objective which tractably enforces that all components are highly informative of reward. Though being task-specific
 892 and compact, RePo is highly sensitive to the hyper-parameters since their objective refer to Lagrangian formulation that
 893 includes various factors. Instead, our framework does not rely on hyper-parameter tuning, where we set all parameters fixed
 894 for all tasks. This further shows the robustness of our framework.
 895
 896

897 D. Experimental Details

898 D.1. Model Architecture Details

900 We have developed our model based on DreamerV3, which employs RSSM to learn state representations and dynamics. We
 901 fix the input image size to 64×64 and use a image encoder which includes a 4-layer CNN with $\{32, 64, 128, 256\}$ channels,
 902 a $(4, 4)$ kernel size, a $(2, 2)$ stride. As a result, our embedding size is 4096.

903 We implement our dynamics model as a hybrid RSSM, which contains an online RSSM for the mask branch and an EMA
 904 RSSM for the raw branch, where the gradients only pass through the online RSSM. The online RSSM is composed of
 905 a GRU and MLPs. The GRU, with 512 recurrent units, is used to predict the current mask recurrent state based on the
 906 previous mask recurrent state, the previous mask posterior stochastic representation, and the previous action. All stochastic
 907 representations are sampled from a vector of softmax distributions, and we use straight-through estimator to backpropagate
 908 gradients through the sampling operation. The EMA RSSM has the same structure as the online RSSM. The size of mask
 909 recurrent states h_t^m is 512 and the size of stochastic representations z_t^m and \hat{z}_t^m is 32×32 . The reward predictor, the
 910 continue predictor, the transition predictor, the value function, and the actor are all MLPs with two hidden layers, each with
 911 512 hidden units. And we use symlog predictions and the discrete regression approach for the reward predictor and the critic.
 912 We use layer normalization and SiLU as the activation function, and update all the parameters using the Adam optimizer.
 913

914 Notably, despite incorporating dual RSSMs, *i.e.*, mask branch and raw branch, in our framework, our model maintains a
 915 slightly smaller overall size (17.54M) compared to the original DreamerV3 (18.22M), which is notable considering the
 916 substantial size of the decoder parameters in DreamerV3. Furthermore, our model is also time-efficient due to the removal
 917 of the time-cost decoder and the use of a shared historical representation for both branches within the framework.
 918

919 D.2. Baselines

920 For DreamerV3, we use an unofficial open-sourced pytorch version of DreamerV3 (NM512, 2023) as the baseline, and we
 921 build our framework on top of it for fair comparison. For RePo and SimSR, we use the official implementation and the
 922 reported hyperparameters in their papers. As for other baselines, we simply adopt the data from results reported in RePo.
 923

924 D.3. Environment Details

926 **DMControl tasks with default settings** This setting consists several continuous control tasks, wherein the agent solely
 927 receives high-dimensional images as inputs. These tasks include *walker_stand*, where a bipedal agent, referred to as “walker”,
 928 is tasked with maintaining an upright position; *walker_walk* and *walker_run*, which require the walker to move forward; and
 929 *cheetah_run*, where a bipedal agent, the “cheetah”, is required to run forward rapidly. We also utilized *cartpole_swingup*, a
 930 task involving a pole and cart system where the goal is to swing up and balance the pole; *reacher_easy*, which involves
 931 controlling a two-link reacher, to reach a target location; and *finger_spin*, where a robotic finger is tasked with continually
 932 rotating a body on an unactuated hinge. We set the time limit to 1000 steps and the action repeat to 2 for all tasks. All
 933 methods were evaluated using 3 different seeds.
 934

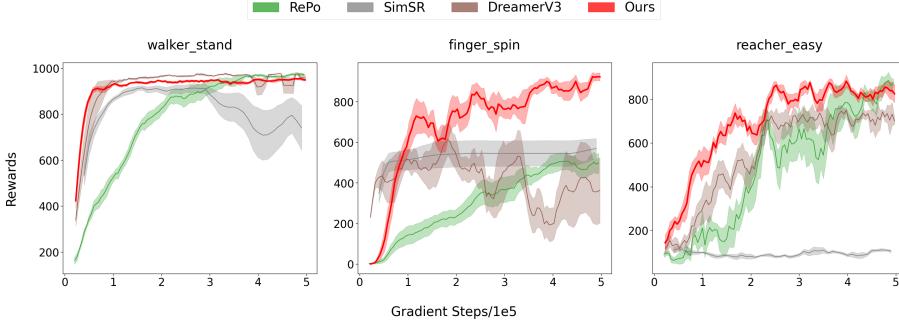


Figure 6. Performance comparison on DMC tasks over 6 seeds in the default setting.

DMControl tasks with distraction settings To evaluate our model’s ability to learn exo-free policy, we test our model in the distraction settings of DMControl. In this setting, we follow DBC (Zhang et al., 2020b) and replace DMControl’s simple static background with 1000 frames grayscale videos from the Kinetics-400 Dataset (Kay et al., 2017), and set the time limit to 1000 steps and the action repeat to 2 for all tasks and evaluate all methods with 3 seeds.

Realistic Maniskill This benchmark is based on the Maniskill 2 (Gu et al., 2023) environment, which encompasses a variety of tasks for the agent to learn to master human-like manipulation skills. To evaluate our model’s ability to learn policy in realistic environments, we follow RePo’s setting and use realistic backgrounds from Matterport (Chang et al., 2017) as distractors. We use 90 scenes from Matterport3D, which are randomly loaded when the environment is reset as distractions for Realistic Maniskill. We set the time limit to 100 steps and the action repeat to 1 for all tasks. All methods were evaluated using 3 different seeds. We test our method and baselines on the tasks *Lift Cube* and *Turn Faucet*: in *Lift Cube*, the agent is required to elevate a cube beyond a specified height, while in *Turn Faucet*, the agent must turn on the faucet by rotating its handle past a target angular distance.

E. Additional Experiments

E.1. Additional performance comparison

We present the learning curve of the methods in the default setting in Figure 6, which is similar to Figure 4 while in a different form of visualization.

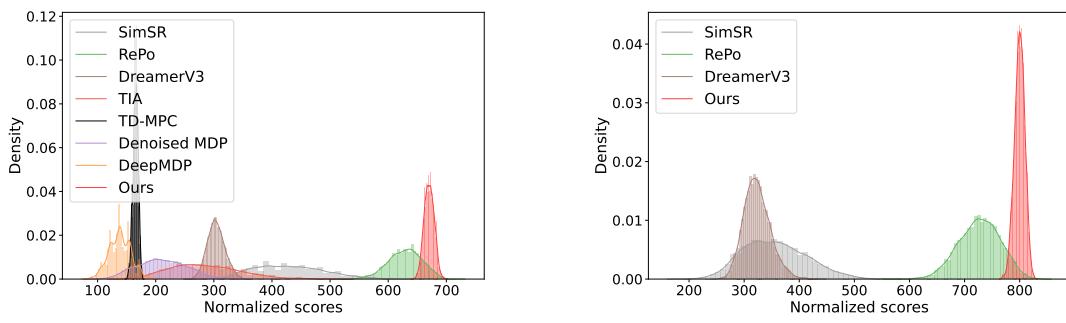


Figure 7. Bootstrapping distributions for uncertainty in IQM (i.e. inter-quartile mean) measurement on DMC tasks in the distraction setting. (**left**) Averaged on 3 seeds. (**right**) Averaged on 6 seeds.

To further statistically illustrate the effectiveness of our model, we present the bootstrapping distributions for uncertainty in IQM (i.e. inter-quartile mean) measurement on DMC tasks in the distraction setting, following from the performance criterion in (Agarwal et al., 2021). Given that the performance results for certain algorithms have been sourced from (Zhu

et al., 2023) and are based on the average across three random seeds, we are unable to calculate the Interquartile Mean (IQM) for all methods with six seeds. Consequently, we present two sets of IQM results in Figure 7. The first set includes all compared methods and is averaged across three seeds. The second set, which we have derived from re-running and evaluating three representative methods, is based on an average across six seeds, providing us with a more robust statistical measure. The result shows that the final performance of our proposed model is statistically better than all other baselines.

E.2. Wall clock time comparison

We compare the wall-clock training time of our method and DreamerV3 in the Realistic Maniskill environment, with the use of a server with NVidia A100SXM4 (40 GB memory) GPU. Figure 8 shows that the running time of our method almost matches DreamerV3, which represents that our model can achieve significant performance improvements at a lower cost, in the presence of exogenous noise, shows that our method can learn effective representations faster.

E.3. Ablation studies

We evaluate the effectiveness of different components of our model by running the ablation experiments on the DMControl’s environment with exogenous noise. All results in this section are averaged across 3 seeds.

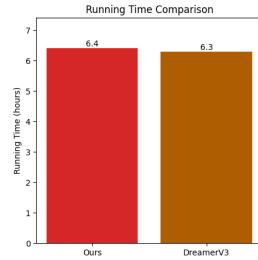


Figure 8. Training Time Comparison on *Lift Cube* task.

Masking-based latent reconstruction and Bisimulation principle Our architecture comprises two main components: masking-based latent reconstruction and a similarity-based objective that follows the bisimulation principle. To assess their effectiveness, we conducted ablation studies by excluding each component individually. Specifically, to evaluate the importance of masking-based latent reconstruction, we removed the mask branch, converting our hybrid RSSM back to a standard RSSM and omitting both the cubic masking and the latent reconstruction loss. To assess the bisimulation principle, we removed only the similarity loss while maintaining all other components.

The results, as shown in Figure 9, reveal that adding just the similarity-based objective to the DreamerV3 framework does not consistently improve sample efficiency across all tasks. This approach often results in the lowest performance, except in the *cartpole_swingup* task. In tasks like *reacher_easy*, the agent fails to develop an acceptable policy, significantly lagging behind in performance compared to other ablations. These findings confirm our theoretical analysis: applying the bisimulation principle directly to model-based agents faces challenges due to the use of an approximate dynamics model for sampling consecutive state representations.

Conversely, utilizing masking-based latent reconstruction generally leads to higher final performance than solely relying on a similarity-based objective. Notably, in nearly half of the tasks, the model with only masking-based latent reconstruction performs comparably to our complete framework, indicating that spatio-temporal information is indeed sparse for these control tasks. Nevertheless, our framework, which includes both components, consistently achieves better performance in most tasks, supporting the necessity of these components. Interestingly, in the *cartpole_swingup* task, the model with only a similarity-based objective outperforms our full framework, suggesting that the integration of both components is not optimal. A possible explanation is that our masking strategy, which is not selectively applied to exogenous noise but rather uses random masking, might inadvertently impact the endogenous state in some contexts.

Normalization for the predictors Our framework incorporates four distinct objectives: latent reconstruction, similarity loss, reward prediction, and episode continuation prediction. For latent reconstruction and similarity loss, we employ normalized state representations because ℓ_2 -normalization ensures that the resulting features are embedded in a unit sphere, which is beneficial for learning state representations. However, the appropriateness of using ℓ_2 -normalization for predicting rewards and episode continuation is not immediately clear. Conventionally, for reward prediction, the exact state representation should be used rather than the normalized one. To investigate this, we conducted an ablation study on the effectiveness of normalization for these two predictors.

The results, illustrated in Figure 10, indicate that normalization may introduce unwanted biases into the predictions, leading to a decrease in performance and increased variance. Therefore, we choose un-normalized representation for reward prediction and continuation prediction.

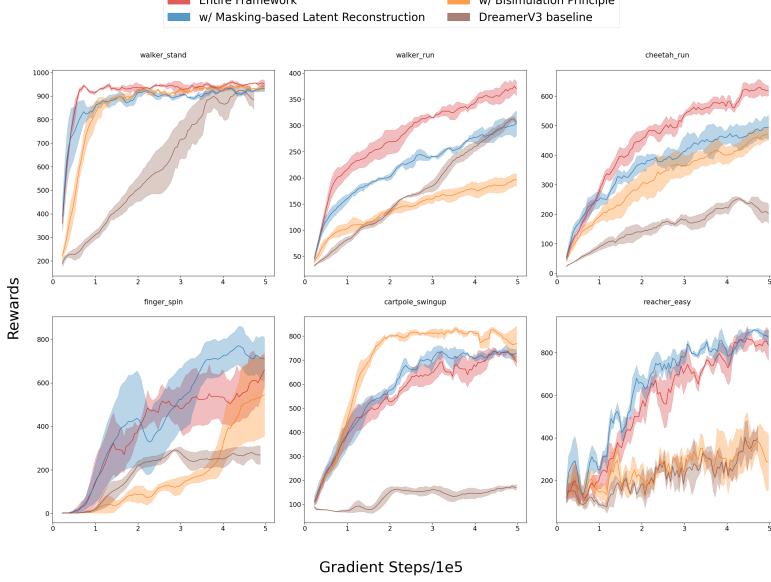


Figure 9. Results of ablation study on masking-based latent reconstruction and the bisimulation principle.

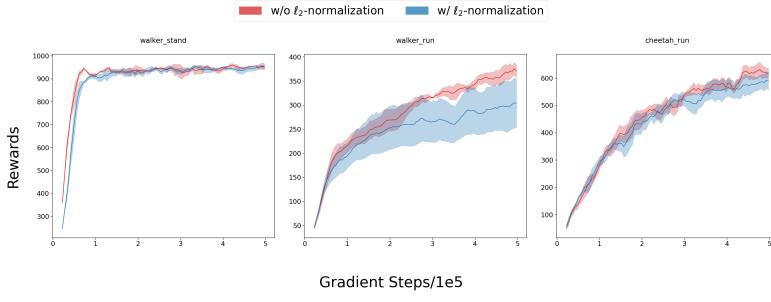


Figure 10. Results of ablation study on ℓ_2 -Normalizion for predictors.

Mask-based Similarity Loss To integrate the masking strategy with similarity loss, we apply a mask to one state in each pair according to the formula, while keeping the other state unmasked. An alternative approach is to use the masked state representation as the current sample pair and the unmasked ones as the consecutive sample pair, *i.e.*,

$$\begin{aligned} \mathcal{L}_{\text{sim}} &:= (d(s_i^m, s_j^m) - \mathcal{F}^\pi d(s_i, s_j))^2 \\ &= \left(d(s_i^m, s_j^m) - \left(|r_{s_i}^\pi - r_{s_j}^\pi| + \gamma d(\hat{s}_{i+1}, \hat{s}_{j+1}) \right) \right)^2 \end{aligned} \quad (16)$$

However, the latter approach may compromise the consistency between the two branches. This is confirmed in Figure 11, which demonstrates that the first approach is more effective, particularly in tasks like *finger_spin*. This effectiveness can likely be attributed to the inherent complexity of the task dynamics. The motion of the manipulated object is influenced not only by the actions of the controllable finger but also by the object's intrinsic inertia, as it undergoes rotational motion. This complexity introduces stochasticity and instability into the environment, posing a significant challenge to dynamics modeling and adversely affecting performance, especially as the policy requires forward-looking dynamics modeling. This ablation study underscores the importance of maintaining consistency between the two branches across various tasks.

Masking ratio We conducted an investigation into the impact of varying mask ratios on the performance of models across different diverse tasks in distraction settings, the result of each task is averaged with three distinct random seeds. The

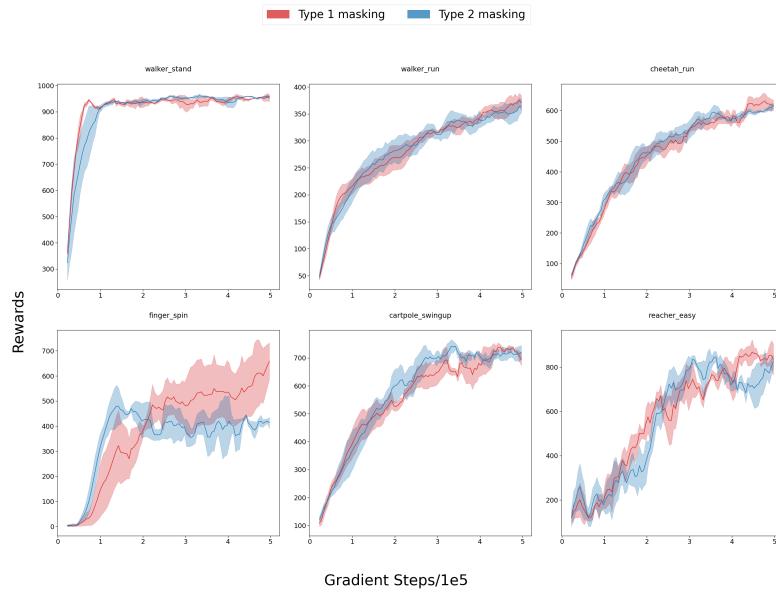


Figure 11. Results of ablation study on masking strategy for similarity loss.

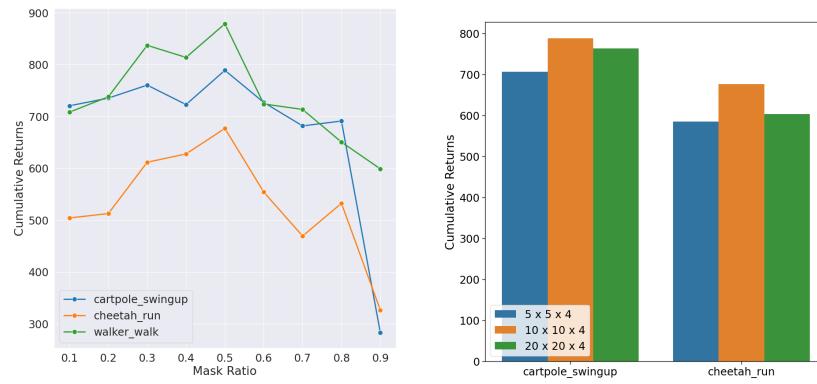


Figure 12. (left) Comparison of different mask ratios in 3 different environments. The final returns are computed at 500k gradient steps updates. (right) Comparison of different patch sizes in 2 different environments. The final returns are computed at 500k gradient steps updates.

masking ratio is selected within a range of 0.1 to 0.9, with the interval of 0.1. The results are depicted in Figure 12. Contrary to the widely-held assumption that image and video data inherently carry a significant degree of superfluous information, our research indicates that the ideal mask ratio for tasks involving sequential control stands at around 0.5. This is notably lower than the nearly 0.9 mask ratio commonly used in computervision domain, as reported in studies such as (He et al., 2022) and (Feichtenhofer et al., 2022). We believe that this discrepancy can be attributed to the control tasks need to retain more spatiotemporal information than CV tasks to facilitate the sequential control tasks.

Cuboid Patch Size We also experimented with different cuboid patch sizes , as $(5 \times 5 \times 4)$, $(10 \times 10 \times 4)$, and $(20 \times 20 \times 4)$ respectively. Throughout the experiments, we maintained a masking ratio of 0.5. The results in Figure 12 indicate that the patch size of $10 \times 10 \times 4$ outperformed both the other two choices. We believe that smaller patch sizes retain unnecessary information, while larger patch sizes may introduce unsuitable masking. Therefore, choosing a moderate patch size is crucial, and in our experiment, we selected $(10 \times 10 \times 4)$ as the default patch size.

E.4. Interpretability visualizations

To verify that our model is indeed capable of filtering task-irrelevant redundancy and learning task-specific features, we utilized the Gradient-weighted Class Activation Mapping (Grad-CAM) technique for feature visualization, as proposed by Selvaraju et al.(Selvaraju et al., 2017). We generate saliency maps for DMC tasks, and then create a binary map, assigning a value of 1 to pixels in the top 5% of intensity values, and 0 otherwise, as illustrated in Figure13.

The results demonstrate that HRSSM can filter out background noise and effectively focuses on the objects that is crucial for control tasks. These results confirm HRSSM’s capability to maintain task-relevant information from visual inputs with exogenous noise within an interpretable manner.

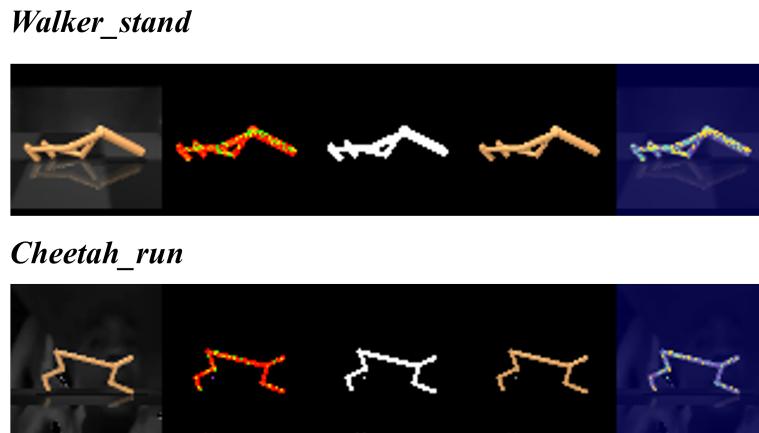


Figure 13. The feature visualization of the learned representations using Grad-CAM.

E.5. More distractions

To evaluate the sample-efficiency and generalization ability of our model, we conduct several different distractions with different nature of noise. Specifically, we have nine different distraction types in total, including the ones we benchmarked in main paper. Examples are in Figure 14. They are:

DMC tasks with default settings This is the default setting of DMC tasks without any distractions. It can be seen as the ideal setting in the realistic tasks.

DMC tasks with distraction settings In this setting, we test the agent in an environment with the background disturbed by the videos from Kinetics dataset (Kay et al., 2017) with the label of driving_car. During the training and evaluation, the environments are both disturbed by the same category of videos, so it is possible that the agent evaluated on the environments that have been seen.

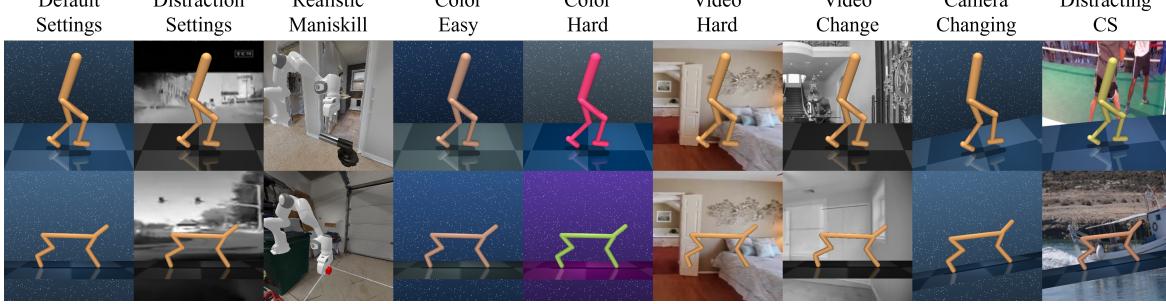


Figure 14. All kinds of distractions.

Table 4. Overview of Distractions.

| Tasks | Distraction |
|-------------------------------------|---|
| DMC tasks with default settings | no distraction |
| DMC tasks with distraction settings | video distraction in background |
| Realistic Maniskill | image distraction in background |
| Color_easy in DMC-GS | colors slightly change for both agent and background |
| Color_hard in DMC-GS | colors dramatically change for both agent and background |
| Video_hard in DMC-GS | video distraction in background, the surface is not visible |
| Video_category_changing | different set of video distraction in background |
| Camera_changing | camera positions change |
| Distracting_CS | all distractions (color, video, camera) |

Realistic Maniskill Similar to DMC tasks with distraction settings. Further, to simulate real-world scenarios, we replace the default background with realistic scenes from the Habitat Matterport dataset (Ramakrishnan et al., 2021), curating 90 different scenes and randomly loading a new scene at the beginning of each episode. So it can be viewed as image distraction in background.

Color_easy in DMC-GS One setting from DeepMind Generalization Benchmark(Hansen & Wang, 2021). We randomize the color of background, floor, and the agent itself, while the colors used are similar to the colors of the original object.

Color_hard in DMC-GS One setting from DeepMind Generalization Benchmark(Hansen & Wang, 2021). Similar to Color_easy, while the colors used is totally different from the colors of the original object.

Video_hard in DMC-GS One setting from DeepMind Generalization Benchmark(Hansen & Wang, 2021). Similar to DMC tasks with distraction settings, while the surface is no longer visible.

Video_category_changing A variation of DMC tasks with distraction settings. During the evaluation, we use a totally different category of videos as background, which makes the testing environments all unseen.

Camera_changing A variation from Distracting Control Suite (Distracting_CS) (Stone et al., 2021) benchmark. We change the span of camera poses and the camera velocity continually throughout an episode.

Distracting_CS Distracting Control Suite (Distracting_CS) (Stone et al., 2021) benchmark is extremely challenging, where camera pose, background, and colors are continually changing throughout an episode. The surface remains visible, such that the agent can orient itself during a changing camera angle.

With the empirical findings presented for the first three distractions in Section 6, our analysis now focuses on the agents' performance against the remaining six distractions. Due to the time limitation, we only have two baseline algorithms tested in these distraction settings in our comparison: SVEA (Hansen et al., 2021) and RePo (Zhu et al., 2023). SVEA is a model-free framework that enhances stability in Q-value estimation by selectively applying data augmentation, optimizing

1265 a modified Bellman equation across augmented and unaugmented data. For RePo, we search several combinations of
 1266 hyperparameters, and choose the best hyperparameter pair in the DMC tasks with distraction setting as default for each task.
 1267 All average returns are averaged by 3 different random seeds. Notably, for Video_category_changing setting, we directly
 1268 evaluate the performance of the agent previously trained on DMC tasks with distraction settings, where we provide the
 1269 agent's final performance metrics rather than a performance progression curve. The results from Figure 15 to Figure 19 and
 1270 Table 5 show that, our model consistently achieve the highest final performance and the best sample-efficiency among the
 1271 most distractions and most tasks, which indicate that our model's robustness and generalization ability across these kinds of
 1272 distractions.
 1273

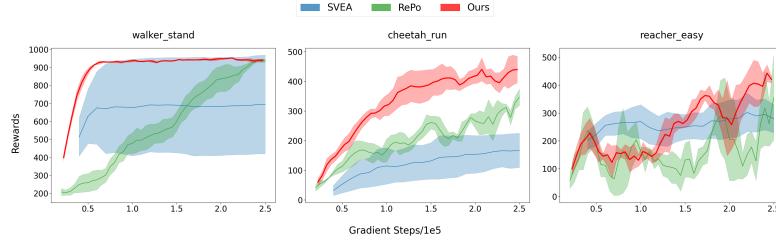


Figure 15. Performance comparison on Color_easy in DMC-GS.

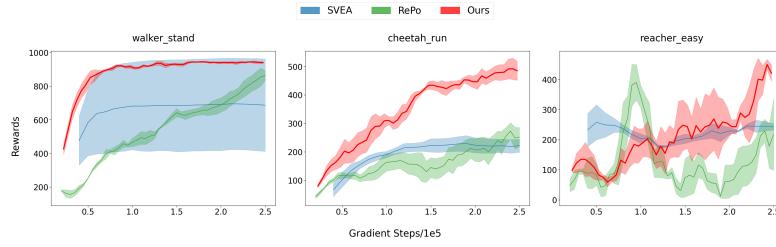


Figure 16. Performance comparison on Color_hard in DMC-GS.

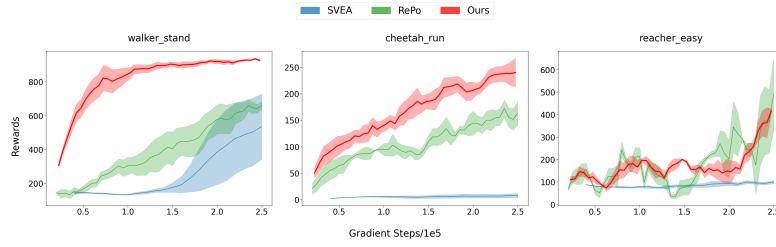


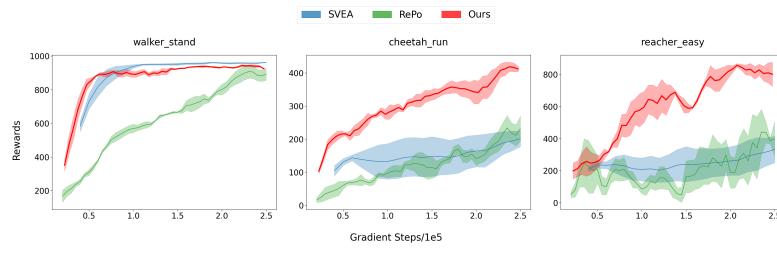
Figure 17. Performance comparison on Video_hard in DMC-GS.

E.6. Analysis of Failure Cases and Bottlenecks of HRSSM

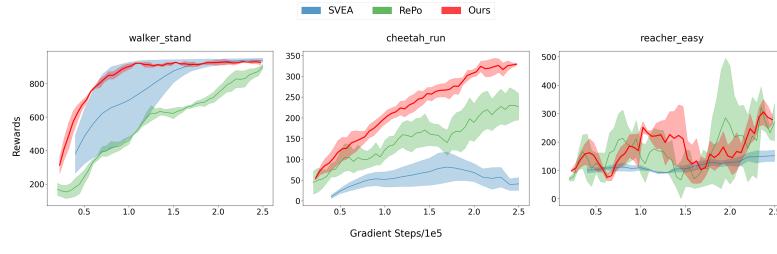
1311 Although our model is effective in most scenarios, as illustrated in previous experiments, there still exist cases that our model
 1312 is not capable of handling well. For instance, the result of *cartpole_swingup* task in our ablation study show that the final
 1313 return of the model that only follows bisimulation principle is higher than our entire model, we consider this can be partially
 1314 attribute to the inappropriate masking. A better masking strategy may help. On the other hand, since our model follows the
 1315 bisimulation principle, it may fail in sparse reward domains, as the fact that the form of bisimulation computation relies on
 1316 bootstrapping with respect to the reward function in recursive terms. We present an evaluation on *ball_in_cup_catch* task in
 1317 Figure 20. The result substantiates the conclusion that our model does not perform well on this kind of tasks. To address this
 1318
 1319

1320
1321
1322
1323
1324 **Table 5.** Final Performance on Video_category_changing. The second column is the final performance on DMC with distraction settings,
1325 where the videos come from the same category. The result shows that our model has good generalization ability to adapt to the unseen
1326 backgrounds.

| Tasks | Rewards on DMC with distraction | Rewards on Video_category_changing |
|------------------|---------------------------------|------------------------------------|
| walker_stand | 946 ± 12 | 963 ± 11 |
| walker_walk | 877 ± 35 | 868 ± 47 |
| walker_run | 390 ± 18 | 406 ± 15 |
| cheetah_run | 652 ± 47 | 628 ± 60 |
| cartpole_swingup | 785 ± 25 | 755 ± 24 |
| reacher_easy | 881 ± 72 | 905 ± 31 |



1335
1336
1337
1338
1339
1340
1341
1342 **Figure 18.** Performance comparison on Camera_changing.



1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360 **Figure 19.** Performance comparison on Distracting_CS.

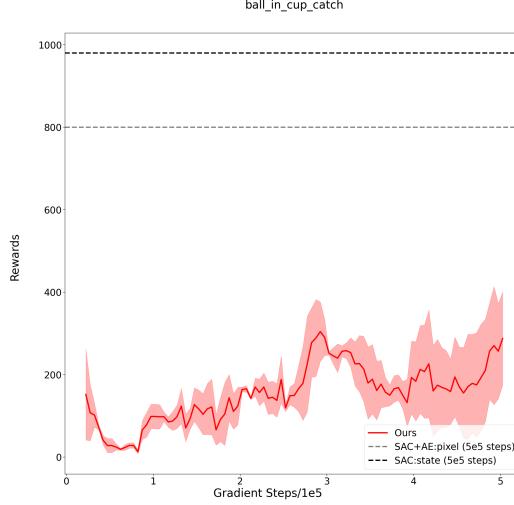


Figure 20. The performance on ball.in.cup.catch task. The gray dotted line is the final performance of SAC+AE (Yarats et al., 2021b) agent with pixel input (same input as ours) that evaluated at 5e5 gradient steps, and the black dotted line is SAC agent with raw state as input that evaluated at 5e5 gradient steps. In this sparse reward task, the performance of our model is not good as the others.

deficiency, employing goal-conditioned RL techniques or implementing reward re-labeling strategies could potentially offer solutions to improve the performance, we leave this to future work.

F. Algorithm

Our training algorithm is shown in Algorithm 1.

1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441 **Algorithm 1** HRSSM
 1442 **Require:** The mask encoder \mathcal{E}_ϕ , the mask posterior model q_ϕ , the mask recurrent model f_ϕ , the mask transition predictor
 1443 $p_\phi(\hat{z}_t^m \mid h_t^m)$, their EMA part $\mathcal{E}'_\phi, q'_\phi, f'_\phi, p'_\phi(\hat{z}_t \mid h_t)$, the reward predictor $p_\phi(\hat{r}_t \mid h_t^m, z_t^m)$ and continue predictor
 1444 $p_\phi(\hat{c}_t \mid h_t^m, z_t^m)$, the critic v_ψ and the actor π_ψ ; the cube masking function $\text{CubeMask}(\cdot)$, the optimizer Optimizer(\cdot, \cdot).
 1445

1446 1: Initialize a replay buffer \mathcal{D} .
 1447 2: Initialize all parameters.
 1448 3: **while** train **do**
 1449 4: **for** update step $c = 1 \dots C$ **do**
 1450 5: // Dynamics learning
 1451 6: Sample B data sequences $\{(a_t, o_t, r_t)\}_{t=k}^{k+T-1}$ from replay buffer \mathcal{D}
 1452 7: Cube masking the observation sequence: $\{o_t^m\}_{t=k}^{k+T-1} \leftarrow \text{CubeMask}(\{o_t\}_{t=k}^{k+T-1})$
 1453 8: Siamese Encoding: $\{e_t^m\}_{t=k}^{k+T-1} \leftarrow \mathcal{E}_\phi(\{o_t^m\}_{t=k}^{k+T-1}), \{e_t\}_{t=k}^{k+T-1} \leftarrow \mathcal{E}'_\phi(\{o_t\}_{t=k}^{k+T-1})$
 1454 9: Compute mask states: $z_t^m \sim q_\phi(z_t^m \mid h_t^m, e_t^m), h_t^m = f_\phi(h_{t-1}^m, z_{t-1}^m, a_{t-1}), \hat{z}_t^m \sim p_\phi(\hat{z}_t^m \mid h_t^m)$
 1455 10: Compute true states: $z_t \sim q'_\phi(z_t \mid h_t^m, e_t), h_t = f'_\phi(h_{t-1}^m, z_{t-1}, a_{t-1}), \hat{z}_t \sim p'_\phi(\hat{z}_t \mid h_t)$
 1456 11: Predict rewards and continuation flags: $\hat{r}_t \sim p_\phi(\hat{r}_t \mid h_t^m, z_t^m), \hat{c}_t \sim p_\phi(\hat{c}_t \mid h_t^m, z_t^m)$
 1457 12: Calculate \mathcal{L}_{dyn} according to Eq. 3
 1458 13: Calculate \mathcal{L}_{rec} according to Eq. 4
 1459 14: Calculate \mathcal{L}_{sim} according to Eq. 5
 1460 15: Calculate $\mathcal{L}_{\text{pred}}$ according to Eq. 6
 1461 16: Calculate total loss $\mathcal{L}(\phi) = \mathbb{E}_{q_\phi} \left[\sum_{t=1}^T (\mathcal{L}_{\text{dyn}}(\phi) + \mathcal{L}_{\text{rec}}(\phi) + \mathcal{L}_{\text{sim}}(\phi) + \mathcal{L}_{\text{pred}}(\phi)) \right]$
 1462 17: Update the encoder's, RSSM's and predictors' parameters: $\mathcal{E}_\phi, q_\phi, f_\phi, p_\phi \leftarrow \text{Optimizer}(\mathcal{E}_\phi, q_\phi, f_\phi, p_\phi, \mathcal{L}(\phi))$
 1463 18: Update the EMA part's parameters: $\mathcal{E}'_\phi \leftarrow m\mathcal{E}_\phi + (1-m)\mathcal{E}'_\phi, q'_\phi \leftarrow mq_\phi + (1-m)q'_\phi, f'_\phi \leftarrow mf_\phi + (1-m)f'_\phi, p'_\phi(\hat{z}_t \mid h_t) \leftarrow mp_\phi(\hat{z}_t \mid h_t) + (1-m)p'_\phi(\hat{z}_t \mid h_t)$
 1464 // Behavior learning
 1465 19: Imagine trajectories $\{(s_t, a_t)\}_{t=k}^{k+H-1}$ from each s_t .
 1466 20: Compute rewards $\{r_t\}_{t=k}^{k+H-1}$, continuation flags $\{c_t\}_{t=k}^{k+H-1}$, values $\{v_t\}_{t=k}^{k+H-1}$ and actions $\{a_t\}_{t=k}^{k+H-1}$
 1467 21: Update the actor π_ψ and the critic v_ψ 's parameters using actor-critic learning.
 1468 22: **end for**
 1469 23: Interact with the environment based on the policy
 1470 24: **end while**

 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484