FOCUS: HIGH-DIMENSIONAL MODEL-BASED RL VIA FOCUSED CONTROL UNIT SAMPLING

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) in high-dimensional continuous state and action spaces often struggles with low learning efficiency and limited exploration scalability. To address this, we introduce FOCUS, a novel model-based RL framework that leverages the insight that effective policies often rely on dynamically focused, sparse control. FOCUS learns preferences over action dimensions to facilitate more targeted and efficient policy learning. It employs a hierarchical decision-making strategy, in which a high-level policy generates binary prompts to activate control units that have more impact on the task performance, while a low-level policy produces actions conditioned on these prompts. To promote behavioral diversity guided by different control-unit preferences, we integrate a diversity-driven objective into the model-based policy optimization process. FOCUS significantly outperforms existing methods on multiple visual control tasks. Furthermore, it facilitates the integration of prior knowledge about the importance of action dimensions, making it particularly effective for complex, high-dimensional tasks.

1 Introduction

Model-based reinforcement learning (MBRL) provides a promising framework for decision-making tasks with high-dimensional observations by modeling environment dynamics (Hafner et al., 2020; 2025). However, its scalability is limited in high-dimensional action spaces, where exploration and policy optimization tend to become inefficient and unstable. In many real-world domains, such as humanoid locomotion, power grid control, and multi-agent collaboration, although the action spaces are high-dimensional, only a sparse subset of dimensions may be crucial for effective decision-making at each step. This insight highlights the need for mechanisms that can selectively attend to key control dimensions during learning, enabling more efficient and focused policy optimization.

A natural approach to mitigating the challenges of high-dimensional action spaces is to exploit their underlying structure. In particular, biasing policy learning toward more informative action dimensions can substantially reduce search complexity and improve efficiency. Several prior methods have explored this direction by introducing hierarchical or factorized action representations. However, many of these approaches rely on *discrete action assumptions* (Saito et al., 2024; Chen et al., 2019; Kumar et al., 2017) or *predefined base-action set* into which the original high-dimensional actions are decomposed (Kim & Dean, 2002; Geißer et al., 2020; Pierrot et al., 2021), thereby limiting their scalability to general continuous control tasks.

To address the under-explored problem of *jointly handling high-dimensional visual observations* and high-dimensional continuous action spaces, we propose FOCUS, a scalable MBRL framework. This setting presents unique challenges, including increased sample complexity and the difficulty of learning effective policies in expansive state-action spaces. Unlike existing MBRL methods (Hafner et al., 2025; Hansen et al., 2023) that treat all action dimensions uniformly, as shown in Figure 1, FOCUS introduces a *control-unit preference learning* mechanism that automatically identifies promising action subspaces and adaptively steers behavior learning toward them, enabling more efficient and focused behavior optimization and exploration.

Specifically, FOCUS represents control-unit preferences using learnable Bernoulli distributions over all action dimensions. The agent presents a *hierarchical decision-making strategy*: (i) A high-level preference policy predicts sampling probabilities and generates binary preference codes

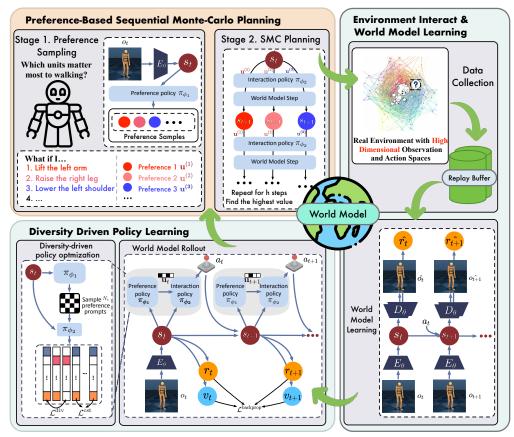


Figure 1: **FOCUS solves high-dimensional continuous control tasks.** The key idea is to improve policy optimization and exploration efficiency via a control-unit preference learning mechanism that adaptively selects and optimizes a dynamic subset of action dimensions.

(the "preference prompt") over all control units; (ii) a low-level interaction policy, conditioned on this prompt, specializes in the selected action subspaces. Both policies are jointly optimized to maximize expected returns over imagined state rollouts produced by the world model. To encourage model-based exploration and alignment between decision layers, we introduce a *diversity-driven policy optimization objective* that maximizes behavioral discrepancies across different control-unit preferences.

Once trained, the agent leverages both policies for real-environment interaction by performing *preference-based Monte Carlo planning*, which reduces possible control-unit sampling bias and significantly improves the decision quality. An additional benefit of our approach is that, when prior knowledge about action dimension priorities is available, it can seamlessly incorporate such inductive biases into the high-level policy.

Experiments on visual control benchmarks with high-dimensional action spaces, *i.e.*, DeepMind Control (Tassa et al., 2018) and MyoSuite (Caggiano et al., 2022), demonstrate that FOCUS outperforms strong model-based and model-free RL approaches in both learning efficiency and final performance.

In summary, the main contributions of this work are as follows:

- We present a unified framework for *high-dimensional MBRL*. It improves existing MBRL methods for efficient learning in large action spaces. Our method dynamically prioritizes promising control units for optimization and planning by modeling the importance of individual action dimensions.
- To enable the interaction policy to respond to varying control-unit preferences, we propose *diversity-driven model-based policy learning*, which leverages a world model to imagine and evaluate trajectories under diverse preference samplings.
- As an additional contribution, we introduce a preference-based sequential Monte Carlo planning
 algorithm, which delivers higher computational efficiency and decision quality than CEM-based
 planning approaches in high-dimensional RL, including TD-MPC2 (Hansen et al., 2023).

2 PROBLEM FORMULATION

Our work targets high-dimensional continuous visual control tasks, where both the observation and action spaces are large and complex, posing significant challenges for efficient learning and exploration. Formally, we formulate the problem as a Partially Observable Markov Decision Process (POMDP), defined by a 5-tuple $(\mathcal{O}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma)$, where: \mathcal{O} is the high-dimensional visual observation space, \mathcal{A} is the high-dimensional continuous action space, \mathcal{R} is the reward function, \mathcal{T} is the state transition dynamics, and γ is the discount factor. In particular, we assume the action space is structured as a Cartesian product of d continuous subspaces:

$$\mathcal{A} = \mathcal{A}^{(1)} \times \mathcal{A}^{(2)} \times \dots \times \mathcal{A}^{(d)}, \quad \text{with } \mathcal{A}^{(i)} \subset \mathbb{R}. \tag{1}$$

An action at time step t is thus represented as a d-dimensional vector $\left[a_t^{(1)}, a_t^{(2)}, \dots, a_t^{(d)}\right]$, where each $a_t^{(i)} \in \mathcal{A}^{(i)}$ is referred to as the i-th sub-action. The objective is to learn a policy $\pi: \mathcal{O} \to \mathcal{A}$, that maximizes the expected discounted cumulative return $\mathbb{E}\left[\sum_{t=1}^T \gamma^{t-1} r_t\right]$.

Challenges. The entanglement of high-dimensional visual observations and high-dimensional continuous actions presents unique and under-explored challenges. Specifically, the state space induced by visual inputs (*e.g.*, raw images) is often noisy, redundant, and partially observable, requiring agents to infer latent dynamics over time. Meanwhile, large continuous action spaces lead to exponentially increasing possibilities for control decisions. Their intersection leads to an expansive and entangled state-action space, making exploration inefficient and value estimation difficult. As a result, existing RL methods often struggle with sample efficiency and generalization in such settings.

3 Method

3.1 OVERVIEW OF FOCUS

To tackle the challenges described above, we propose an MBRL method guided by learned controlunit preferences. The overall training and environment interaction pipeline consists of three stages:

- (i) World model learning: We first train a world model to predict latent state transitions and reward signals from observation-action-reward tuples. We use DreamerV3 (Hafner et al., 2025) as the backbone for world modeling. Detailed model specifications are provided in the appendix.
- (ii) Diversity-driven model-based policy learning: As shown in Algorithm 1, over future trajectory rollouts (i.e., latent state imaginations) generated by the world model, FOCUS constructs a hierarchical policy learning architecture. At its core is a high-level decision-making module that outputs a preference distribution over all control units¹. At each imagination step, the high-level policy samples a binary control-unit preference from Bernoulli distributions (termed a "preference prompt"). This prompt selects a sparse subset of control units, guiding the low-level policy network to generate specific behaviors within the original high-dimensional action space. The high-level preference policy and the low-level interaction policy are jointly optimized to maximize expected returns over the imagination horizon. Furthermore, we encourage behavioral diversity across different preference prompts, promoting more effective exploration and stronger alignment between hierarchical decisions.
- (iii) *Preference-based Monte Carlo planning:* During environment interaction, the agent samples multiple candidate preference prompts. For each prompt, it performs multi-step Monte Carlo planning to generate potential future trajectories. The final action is then selected from these candidates based on trajectory weights computed from their cumulative advantage estimates.

3.2 CONTROL-UNIT PREFERENCE SAMPLING

Our approach is motivated by the observation that high-dimensional control tasks often exhibit inherent sparsity—at any given time, some control units more significantly contribute to task performance than others. For example, in the HumanoidStand task, the agent begins by pushing itself upward

¹Each *control unit* corresponds to a subset of action dimensions for a specific functional component (*e.g.*, a joint in a humanoid robot), potentially spanning multiple dimensions to reflect its degrees of freedom.

Algorithm 1: Full training algorithm of FOCUS, with Diversity-driven model-based policy learning, and Preference-based SMC Planning

```
164
          1 Initialize: world model \theta, hierarchical actor \{\phi_1, \phi_2\}; critic \psi, replay buffer \mathcal{B} with random episodes.
165
          2 while not converged do
166
                    for update step c = 1 \dots C do
167
                           // World model learning
          4
                           Draw data sequences \{(o_t, a_t, r_t)\}_{t=1}^T \sim \mathcal{B}.
          5
169
                           Update the dynamics model p_{\theta} and reward model r_{\theta} using DreamerV3 world model objectives.
                             / Diversity-driven model-based policy learning
170
                          Draw a random observation o_1 \sim \mathcal{B}.
171
                          for time step t = 1 \dots H - 1 do
172
                                 Sample N_c control-unit preference promts \{\mathbf{u}_{t,n}\}_{n=1}^{N_c} \sim \pi_{\phi_1}(s_t).
          10
173
                                 Generate \{a_{t,n} \sim \pi_{\phi_2}(s_t, \mathbf{u}_{t,n})\}_{n=1}^{N_c}.
          11
174
                                 Predict \{s_{t+1,n} \sim p_{\theta}(s_t, a_{t,n})\}_{n=1}^{N_c} and \{r_{t,n} \sim r_{\theta}(s_{t+1,n})\}_{n=1}^{N_c}.
Compute \mathcal{L}_t^{\text{div}} and \mathcal{L}_t^{\text{cst}} over \{a_{t,n}\}_{n=1}^{N_c} according to Eq. equation 3 and Eq. equation 4.
          12
175
          13
176
                                 Select s_{t+1,n} with the highest v_{\psi}(s_{t+1,n}) as s_{t+1}.
          14
177
          15
                          Update \pi_{\phi_1}, \pi_{\phi_2} and v_{\psi} using Eq. equation 5.
          16
179
          17
                    // Environment interaction with sequential Monte Carlo planning
         18
                    o_1 \leftarrow \text{env.reset}()
          19
181
                     \  \, {\bf for} \  \, {\it time} \  \, {\it step} \  \, t=1\dots T \  \, {\bf do} \\
         20
                          Compute posterior state s_t \sim q_\theta(o_t).
          21
183
                          Sample N_p preference promts \{\mathbf{u}_n\}_{n=1}^{N_p} \sim \pi_{\phi_1}(\cdot \mid s_t).
Initialize N_p planning trajectories \{\tau_n\}_{n=1}^{N_p} starting from \{\hat{s}_{t,n}\}_{n=1}^{N_p} = s_t.
          22
          23
185
                           Initialize importance weights for planning trajectories \{w_n = 1\}_{n=1}^{N_p}.
          24
186
                          for planning steps h = 0 \dots L - 1 do
187
                                Generate \{\hat{a}_{t+h,n} \sim \pi_{\phi_2}(\hat{s}_{t+h,n},\mathbf{u}_n)\}_{n=1}^{N_p}.
Predict \{\hat{s}_{t+h+1,n} \sim p_{\theta}(\hat{s}_{t+h,n},\hat{a}_{t+h,n})\}_{n=1}^{N_p} and \{\hat{r}_{t+h,n} \sim p_{\theta}(\hat{s}_{t+h+1,n})\}_{n=1}^{N_p}.
188
          27
189
                                 Compute step-wise advantages \{A_{t+h,n} = \hat{r}_{t+h,n} + \gamma v_{\psi}(\hat{s}_{t+h+1,n}) - v_{\psi}(\hat{s}_{t+h,n})\}_{n=1}^{N_p}.
190
                                 Update weights \{w_{t+h+1,n} = w_{t+h,n} \cdot \exp(A_{t+h,n})\}_{n=1}^{N_p}.
          29
192
                           Sample \tau_n from a categorical distribution based on weights \{w_{t+L,n}\}_{n=1}^{N_p}.
193
                          Select the first action \hat{a}_{t,n} in \tau_n as a_t.
          32
194
                          r_t, o_{t+1} \leftarrow \text{env.step}(a_t)
          33
         34
196
                    Add experience to the replay buffer \mathcal{B} \leftarrow \mathcal{B} \cup \{(o_t, a_t, r_t)_{t=1}^T\}.
          35
             end
```

from the ground using support from its legs and arms. In such moments, joints such as the ankles and wrists play a dominant role, while many other action dimensions contribute minimally. This observation suggests that exploration can be made more efficient by selectively attending to a small subset of all control units. Accordingly, we introduce a high-level policy π_{ϕ_1} that takes latent states as inputs and generates a binary prompt $\mathbf{u}_t \in \mathbb{R}^d$, where $\mathbf{u}_t^{(i)} = 1$ indicates that the i-th action dimension is currently important and should be prioritized. This design leads to a hierarchical policy structure, where the target policy π_{ϕ_2} is conditioned on both the latent state and the preference prompt:

Preference policy:
$$\mathbf{u}_t \sim \pi_{\phi_1}(s_t)$$
, Interaction policy: $a_t \sim \pi_{\phi_2}(\mathbf{u}_t, s_t)$, (2)

where \mathbf{u}_t is sampled from d independent Bernoulli distributions whose parameters are jointly predicted by the high-level policy π_{ϕ_1} , enabling the modeling of correlations across all control units. Intuitively, we expect the target policy π_{ϕ_2} to adapt its behavior based on the high-level guidance \mathbf{u}_t , enabling more focused and effective exploration within the subspace of prioritized action dimensions.

3.3 DIVERSITY-DRIVEN MODEL-BASED POLICY LEARNING

To promote an efficient exploration of potential policies guided by diverse control-unit activations, we introduce a focused policy optimization approach within a model-based behavior learning framework, as shown in Figure 1. We first learn a world model parameterized by θ , which includes an encoder $p_{\theta}(s_t \mid o_t)$, a latent transition model $p_{\theta}(s_{t+1} \mid s_t, a_t)$, and a reward predictor $r_{\theta}(s_t)$. Given an initial latent state encoded from a random observation in the replay buffer, FOCUS samples a preference prompt \mathbf{u}_t , and generates an preference-conditioned action a_t based on the current rollout state \hat{s}_t . The action a_t is then used to update the imagination trajectory via the world model. The hierarchical policy $\pi_{\phi_{1,2}}$ and the value model v_{ψ} are optimized over imagined trajectories of length H.

However, directly maximizing the expected value function, as in Dreamer-style objectives, is insufficient for inducing low-level behaviors that are responsive to high-level preferences. Without further regularization, the actor may fail to interpret and follow the high-level guidance effectively. To address this, we design a diversity-driven learning objective. The key idea is that the low-level actor should produce diverse behaviors for control units selected by different preference prompts, while maintaining consistent behaviors for those not selected. Specifically, we sample N_c preference prompts $\{\mathbf{u}_{t,n}\}_{n=1}^{N_c}$ from $\pi_{\phi_1}(s_t)$, and and generate corresponding action vectors $\{a_{t,n}\}_{n=1}^{N_c}$. Each action is decomposed as $a_{t,n} = [a_{t,n}^{(1)}, \ldots, a_{t,n}^{(d)}]$ along d dimensions. We first define a diversity-promoting loss to encourage the low-level policy to respond distinctly to different preference selections:

$$\mathcal{L}_{t}^{\text{div}} = \sum_{i=1}^{d} \left[-\sum_{n=1}^{N_{c}} \mathbb{I}(\mathbf{u}_{n}^{(i)} = 1) \sum_{k=1}^{N_{c}} \mathbb{I}(\mathbf{u}_{k}^{(i)} = 0) \cdot \text{KL} \left[\pi(a_{n}^{(i)} \mid s, \mathbf{u}_{n}) \parallel \text{sg}(\pi(a_{k}^{(i)} \mid s, \mathbf{u}_{k})) \right] \right], (3)$$

where $\mathbb{I}(\cdot)$ is the indicator function, $\operatorname{sg}(\cdot)$ denotes stop-gradient, and we ignore the time step index for clarity. Meanwhile, we introduce a consistency loss to penalize divergence among actions on the same dimensions not selected by the preference prompt:

$$\mathcal{L}_{t}^{\text{cst}} = \sum_{j=1}^{d} \left[\sum_{n=1}^{N_{c}} \mathbb{I}(\mathbf{u}_{n}^{(j)} = 0) \sum_{k=1}^{N_{c}} \mathbb{I}(\mathbf{u}_{k}^{(j)} = 0) \cdot \text{KL} \left[p(a_{n}^{(j)} \mid s, \mathbf{u}_{n}) \parallel \text{sg}(p(a_{k}^{(j)} \mid s, \mathbf{u}_{k}))] \right] \right]. \quad (4)$$

Together, these losses impose a contrastive structure in the high-dimensional action space, enhancing exploration and ensuring better alignment between low- and high-level policies. Importantly, we do not stop the gradient at \mathbf{u}_n in the above regularization terms. Instead, we allow π_{ϕ_1} to be co-optimized with π_{ϕ_2} . This design further promotes exploration, as $\mathcal{L}_t^{\text{div}}$ rewards diverse control-unit activations and prevents the high-level policy from collapsing to identical but biased activation vectors.

For the sake of computational efficiency, as shown in Alogrithm 1, we choose $s_{t+1,n} \sim p_{\theta}(s_t, a_{t,n})$ with the highest predicted value $v_{\psi}(s_{t+1,n})$ as s_{t+1} among multiple candidates at each imagination step. In practice, we compute the diversity-driven objectives only at the first step of imagination. Despite this simplification, we observe a significant performance gain over versions without diversity-driven regularization, while incurring almost no additional training overhead.

Finally, the full objective is defined as follows, where the actor uses value backpropagation through dynamics as in DreamerV3 for better continuous control, and the critic uses a λ -return target:

Actor:
$$\mathcal{L}(\phi_1, \phi_2) = \mathbb{E}_{p_{\psi}, p_{\theta}} \left[\sum_{t=1}^{H-1} \left[\mathcal{L}_t^{\text{div}} + \mathcal{L}_t^{\text{cst}} - V_t^{\lambda} - \eta \mathcal{H} \left(a_t \mid \hat{s}_t, \text{sg}(\mathbf{u}_t) \right) \right] \right],$$

$$\text{Critic:} \quad \mathcal{L}(\psi) = \mathbb{E}_{p_{\phi_1}, p_{\phi_2}, p_{\theta}} \left[\sum_{t=1}^{H-1} \frac{1}{2} (v_{\psi}(\hat{s}_t) - \text{sg}(V_t^{\lambda}))^2 \right],$$
(5)

where the entropy term $\mathcal{H}(\cdot)$ further encourages exploration conditioned on preference prompts. Notably, both levels of policies are jointly optimized towards maximizing the imagined values.

3.4 Preference-Based Monte Carlo Planning

The hierarchical policy involves two sources of stochasticity—from both the high-level preference sampling and the low-level action generation—which may increase variance during real-environment

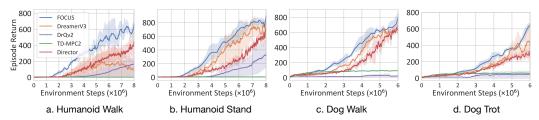


Figure 2: **Results on DeepMind Control Suite.** All results are averaged over three training seeds.

interaction. To mitigate this, we develop an online planning scheme that efficiently reasons over multiple future state-action trajectories, each conditioned on a different control-unit preference prompt. We term this method *preference-based Monte Carlo planning*, a modified form of sequential Monte Carlo planning (Piché et al., 2018) designed for high-dimensional action spaces.

As shown in Algorithm 1, at each interaction step, we sample N_p preference prompts $\{\mathbf{u}_n\}_{n=1}^{N_p}$ based on the current state s_t , and initialize N_p planning trajectories with uniform weights $w_n=1$. At each planning step t+h, we reuse \mathbf{u}_n to sample low-level actions $\hat{a}_{t+h,n}\sim\pi_{\phi_2}(\hat{s}_{t+h,n},\mathbf{u}_n)$. This contrasts with model-based policy learning, where \mathbf{u}_n is typically re-sampled at each imagination step. Instead, we preserve temporal consistency in the sampled control-unit prompts during planning. We then use the learned dynamics p_θ to rollout future states $s_{t+h:t+L,n}$. We evaluate these candidate trajectories via advantage estimation and update their weights accordingly:

$$A_{t+h,n} = \hat{r}_{t+h,n} + \gamma v_{\psi}(\hat{s}_{t+h+1,n}) - v_{\psi}(\hat{s}_{t+h,n}), \quad w_{t+h+1,n} = w_{t+h,n} \cdot \exp(A_{t+h,n}).$$
 (6)

At the end of the planning horizon, we sample one trajectory from the categorical distribution defined by $\{w_{t+L,n}\}_{n=1}^{N_p}$ and execute its first action in the environment.

Empirically, this approach significantly outperforms prior sample-based planning methods (*e.g.*, MPPI in TD-MPC2 (Hansen et al., 2023)) in both computational efficiency and decision quality. Unlike MPPI, our method avoids online policy refitting, making it more suitable for high-dimensional control tasks.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

We conduct experiments on two RL environments with visual inputs, namely DeepMind Control Suite (Tassa et al., 2018) and MyoSuite (Caggiano et al., 2022). We present empirical results on seven tasks from these environments with high-dimensional continuous action spaces. Detailed information on the selected tasks with their action space sizes is provided in Appendix A.

We compare FOCUS with: (1) *DreamerV3* (Hafner et al., 2025), a strong model-based RL method in high-dimensional observation space; (2) *TD-MPC2* (Hansen et al., 2023), a decoder-free model-based RL baseline that employs the CEM-based planning method. (3) *DrQv2* (Yarats et al., 2021a), a data augmentation-based visual RL approach; and (4) *Director* (Hafner et al., 2022), another hierarchical MBRL method that conditions low-level policy on goals generated by the high-level autoencoder. All models are trained under a fixed environment interaction budget. For two tasks in the Humanoid suite, each model is trained for 8M frames, and for the Dog suite, 6M frames. Tasks in MyoSuite allow interactions for 1M frames. Average episode returns are computed over 10 episodes per seed.

4.2 Domain #1: DeepMind Control Suite

As shown in Figure 2, FOCUS consistently outperforms baseline methods across all four tasks. Notably, it achieves a substantial improvement over DreamerV3 in HumanoidWalk. In both Humanoid tasks, the learning curves show a prolonged warm-up phase before meaningful improvement, reflecting the challenge of locating sparse high-reward regions within large decision spaces. Compared with **DreamerV3**, our model exhibits an earlier performance rise, demonstrating its ability to reduce sample complexity and efficiently explore targeted reward regions. In Dog tasks, rewards are less sparse than in Humanoid tasks, and DreamerV3 is already capable of sampling high-reward

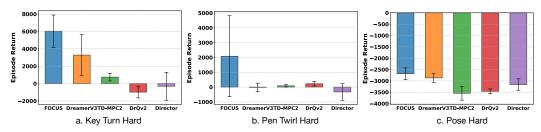


Figure 3: **Results on the MyoSuite hard tasks.** Episode returns are averaged over 10 episodes per seed, for a total of 30 runs. FOCUS outperforms all baselines across every task.

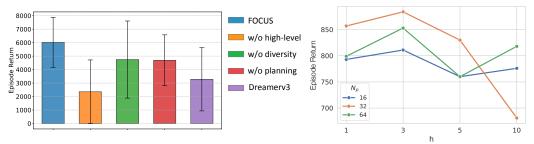


Figure 4: Ablation studies. Left: Effect of each model component (Task: KeyTurnHard). Right: Effect of planning horizon h and candidate trajectories N_p in our planner (Task: HumanoidWalk).

regions. This is reflected in its steadily increasing learning curves, suggesting that despite the larger number of action dimensions, the high-reward regions in the action space are sufficiently large for agents with random policies to sample with reasonable probability. In this setting, exploration is not the primary bottleneck, which limits potential performance gains. For **TD-MPC2**, it fails to achieve meaningful progress on humanoid tasks, highlighting its difficulty in effectively exploring the vast search space, particularly under high-dimensional visual input. The comparison with **Director** highlights the impact of the diversity-driven objective: while both methods employ hierarchical policies, FOCUS explicitly prioritizes subsets of the action space, resulting in higher exploration efficiency. We showcase the policy evaluation of FOCUS, DreamerV3 and Director in Figure 5a, with full results in Appendix B.1. As illustrated, the policy obtained by FOCUS adjusts humanoid's posture in fewer time steps to start walking (t = 50), meanwhile maintaining its control more stable than DreamerV3 which falls to the ground (t = 200) after standing up. Additionally, DrQv2 requires roughly 4M environment steps before its learning curve begins to rise, suggesting that augmenting failing episodes early has little effect on sample efficiency.

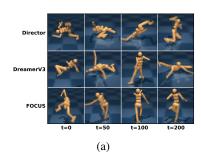
Further results. We compare the planning efficiency of FOCUS with TD-MPC2 in Appendix B.2, and assess its generalizability on tasks with low-dimensional action spaces in Appendix B.4.

4.3 Domain #2: MyoSuite

Performance results on MyoSuite are shown in Figure 3. FOCUS outperforms all baselines on KeyTurnHard and PenTwirlHard. On KeyTurnHard, it improves over DreamerV3, demonstrating higher efficiency in exploring high-reward regions within the same environment step budget On PenTwirlHard, FOCUS successfully discovers sharply peaked reward regions, whereas other baselines struggle in flatter areas; DreamerV3 fails to solve this task. The integration of the diversity-promoting loss \mathcal{L}_t^{div} and consistency loss \mathcal{L}_t^{cst} systematically biases sampling toward high-reward regions, enabling efficient policy learning and sustained performance gains All methods face challenges on PoseHard, though FOCUS achieves slightly better results, thanks in part to \mathcal{L}_t^{cst} preventing overfitting as discussed in Section 3.3.

4.4 ABLATION STUDIES AND HYPERPARAMETER ANALYSES

We validate the proposed modules on KeyTurnHard in MyoSuite, considering three variants of our method: (1) FOCUS without π_{ϕ_1} , reducing it to a DreamerV3 agent planning in the original high-dimensional action space; (2) FOCUS without \mathcal{L}_t^{div} and \mathcal{L}_t^{cst} ; and (3) FOCUS without planning. These variants are compared with the full FOCUS and DreamerV3 to assess the contribution of each



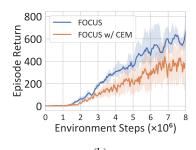


Figure 5: (a) Visualization of the learned policy for HumanoidWalk. FOCUS features more efficient and stable policy learning. (b) Ablation on substituting CEM for SMC planning in HumanoidWalk. Our design demonstrates more efficiency in exploring high-dimensional action space.

component. As shown in Figure 4, FOCUS consistently outperforms all other variants, demonstrating that each component contributes to its success. The variant without the high-level policy π_{ϕ_1} is equivalent to DreamerV3 with the addition of a Sequential Monte Carlo planner during online interaction, yet it performs worse than DreamerV3. We attribute this to the vulnerability of direct planning in the original high-dimensional action space, whereas FOCUS plans in the high-level space, providing more effective guidance. Furthermore, both the diversity loss and the planning algorithm are crucial contributors to FOCUS's performance. We also examine different variants across planning horizons h and number of candidate trajectories N_p , highlighting the h=3 and $N_p=32$ setting. Although increasing N_p substantially raises the computational cost, it does not guarantee consistent improvements in policy performance. Additionally, we compare FOCUS to its Cross Entropy Method(CEM)-based variant in Figure 5b. As illustrated, SMC enables faster and more effective learning compared to CEM. we argue that prior sample-based planning methods like CEM which requires online refitting of policy distribution, are not sufficiently capable of modeling target distributions in a high-dimensional action space and thus limiting its effectiveness.

4.5 VISUALIZATION OF LEARNED CONTROL-UNIT PREFERENCES

We visualize the preference prompts sampled on HumanoidStand in Figure 6, using the model trained immediately after the rising point (*i.e.*, after 3M environment steps) and recording preference prompts over a single episode. This visualization serves two purposes: (1) to illustrate the temporal evolution of control preferences as the agent attempts to stand, and (2) to qualitatively assess whether the learned preference policy captures meaningful, structured sparsity across action dimensions. From the sampled preference prompts, we observe distinct patterns in certain control units that are consistently activated during critical stabilization moments. These results support the hypothesis that FOCUS exploits structured sparsity to focus exploration on behaviorally relevant subspaces, aligning with our goal of reducing exploration overhead in high-dimensional control tasks.

5 RELATED WORK

Model-based RL. Model-based RL (Ha & Schmidhuber, 2018; Hafner et al., 2025; Hansen et al., 2023; Moerland et al., 2023) shows great potential in addressing the data efficiency issue compared to model-free methods (Yarats et al., 2021b; Kostrikov et al., 2021; Laskin et al., 2020), particularly in online learning. By leveraging a learned world model, these methods enable efficient policy optimization through planning (Nguyen et al., 2021; Zhao et al., 2021; Hansen et al., 2023) or Qlearning (Hafner et al., 2020; Wang et al., 2022). However, in high-dimensional observation and action spaces, world models struggle to identify sparse high-reward regions due to their uniform treatment of all action dimensions. This limitation burdens the learning process, as undirected exploration in vast action spaces leads to inefficiency. FOCUS addresses this challenge by adaptively selecting a control unit, thereby guiding policy exploration and optimization within a reduced subspace.

Hierarchical RL. Prior studies in hierarchical structured action space (Kumar et al., 2017; Chen et al., 2019; Tang & Agrawal, 2020; Saito et al., 2024) usually assume a multi-stage policy that divides the searching space into smaller parts. However, this structure is mostly restricted to tasks with discrete action spaces and requires additional discretization when confronted with continuous ones. Hierarchical policy learning in behavior learning (Pateria et al., 2021) has been used to enable efficient

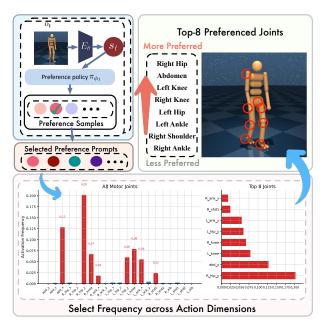


Figure 6: **Visualization of learnt preference prompts.** These results suggest that FOCUS effectively leverages structured sparsity to focus exploration on task-relevant subspaces, aligning with our goal of reducing exploration overhead in high-dimensional control.

long-horizon reasoning by generating temporal abstractions through high-level policy (Vezhnevets et al., 2017; Gumbsch et al., 2024; Gürtler et al., 2021). Director (Hafner et al., 2022) is a model-based method featuring a hierarchical policy structure that generates a goal with a high-level autoencoder first. Yet in our work, we have made use of the hierarchical structure to model the sparsity nature of the high-dimensional action space, thus reducing the searching space to a single control unit and guiding the exploration in a high-dimensional continuous action space.

RL with high-dimensional action space. Despite the hierarchical structure mentioned above, similar structure-exploiting methods have been proposed, faced with challenges in high-dimensional action space, among which RL in factored action space (Osband & Van Roy, 2014) serves as a favorable direction. These methods make use of the loose-coupling action space and decompose the original MDP into several sub-questions and prove especially efficient in multi-agent settings (Guestrin et al., 2001; Mahajan et al., 2021; Peng et al., 2021). The limitation lies in the strong assumption that access to a pre-defined base-action set into which the original high-dimensional action could be decomposed is granted, demonstrating limited efficacy in single-object control tasks such as humanoid. Additionally, learning-based structure discovery methods (Tavakoli et al., 2017; 2018; Van de Wiele et al., 2020) are similar in spirit to our approach, as they model each action dimension with separate actors. However, this design introduces computational overhead and scalability challenges, and often ignores interdependencies between action dimensions, which can lead to instability during training.

6 Conclusions and Limitations

In this paper, we introduced FOCUS, a novel approach for high-dimensional model-based reinforcement learning (MBRL) that addresses the challenges of learning and exploration efficiency in large action spaces. FOCUS leverages a hierarchical policy structure, where a high-level policy learns to generate sparse control-unit activations, identifying key action dimensions at each decision point; while a low-level policy operates within these preferred subspaces to optimize behavior using imagined rollouts from a world model. Our empirical results showed that FOCUS consistently outperforms baselines on high-dimensional control tasks while maintaining high learning efficiency.

Despite its strong performance, FOCUS also has certain limitations. First, the effectiveness of the learned preference prompts closely tied to the quality of the world model, as inaccuracies in imagined rollouts can lead to suboptimal high-level guidance. Second, the hierarchical decision-making mechanism introduces extra parameters and training complexity. Exploring effective preference representations can be a promising direction for future work.

ETHICS STATEMENT

This work proposes an RL framework focusing on dealing with high-dimensional action spaces. The research does not involve human subjects, sensitive personal data, or any newly collected datasets. All experiments were conducted within simulated environments, ensuring that no privacy or security concerns are introduced.

Our study does not directly present harmful applications. The contributions are methodological and algorithmic in nature, aiming to deepen the understanding of efficient policy learning in high-dimensional spaces. We have carefully read and followed the Code of Ethics, and are committed to principles of research integrity. No conflicts of interest or external sponsorship had influence on this work.

REPRODUCIBILITY STATEMENT

This work proposes a novel model, FOCUS, to deal with RL in tasks with high-dimensional observation and action input. To facilitate reproducibility, we provide the complete implementation code as part of the supplementary materials. Detailed descriptions of the model is included in Section 3.2, Section 3.3 and the planning algorithm is included in Section 3.4. We also describe the detailed experimental settings in Section 4.1 and Appendix A.

REFERENCES

- Vittorio Caggiano, Huawei Wang, Guillaume Durandau, Massimo Sartori, and Vikash Kumar. Myosuite–a contact-rich simulation suite for musculoskeletal motor control. *arXiv preprint arXiv:2205.13600*, 2022.
- Haokun Chen, Xinyi Dai, Han Cai, Weinan Zhang, Xuejian Wang, Ruiming Tang, Yuzhou Zhang, and Yong Yu. Large-scale interactive recommendation with tree-structured policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 3312–3320, 2019.
- Florian Geißer, David Speck, and Thomas Keller. Trial-based heuristic tree search for mdps with factored action spaces. In *Proceedings of the International Symposium on Combinatorial Search*, volume 11, pp. 38–47, 2020.
- Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. *Advances in neural information processing systems*, 14, 2001.
- Christian Gumbsch, Noor Sajid, Georg Martius, and Martin V. Butz. Learning hierarchical world models with adaptive temporal abstractions from discrete latent dynamics. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=TjCDNssXKU.
- Nico Gürtler, Dieter Büchler, and Georg Martius. Hierarchical reinforcement learning with timed subgoals. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 21732–21743. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/b59c21a078fde074a6750e91ed19fb21-Paper.pdf.
- David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Proceedings* of the 32nd International Conference on Neural Information Processing Systems, NIPS'18, pp. 2455–2467, Red Hook, NY, USA, 2018. Curran Associates Inc.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *ICLR*, 2020.
- Danijar Hafner, Kuang-Huei Lee, Ian Fischer, and Pieter Abbeel. Deep hierarchical planning from pixels. *Advances in Neural Information Processing Systems*, 35:26091–26104, 2022.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *Nature*, 640:647–653, 2025.

- Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for continuous control. *arXiv preprint arXiv:2310.16828*, 2023.
- Kee-Eung Kim and Thomas Dean. Solving factored mdps with large action space using algebraic decision diagrams. In *Pacific Rim International Conference on Artificial Intelligence*, pp. 80–89. Springer, 2002.
 - Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *ICLR*, 2021.
 - A Kumar, K Swersky, and G Hinton. Feudal learning for large discrete action spaces with recursive substructure. In *Proceedings of the NIPS Workshop Hierarchical Reinforcement Learning, Long Beach, CA, USA*, volume 9, 2017.
 - Michael Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. In *NeurIPS*, 2020.
 - Anuj Mahajan, Mikayel Samvelyan, Lei Mao, Viktor Makoviychuk, Animesh Garg, Jean Kossaifi, Shimon Whiteson, Yuke Zhu, and Animashree Anandkumar. Tesseract: Tensorised actors for multi-agent reinforcement learning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 7301–7312. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/mahajan21a.html.
 - Thomas M Moerland, Joost Broekens, Aske Plaat, Catholijn M Jonker, et al. Model-based reinforcement learning: A survey. *Foundations and Trends*® *in Machine Learning*, 16(1):1–118, 2023.
 - Tung D Nguyen, Rui Shu, Tuan Pham, Hung Bui, and Stefano Ermon. Temporal predictive coding for model-based planning in latent space. In *International Conference on Machine Learning*, pp. 8130–8139. PMLR, 2021.
 - Ian Osband and Benjamin Van Roy. Near-optimal Reinforcement Learning in Factored MDPs. *arXiv e-prints*, art. arXiv:1403.3741, March 2014. doi: 10.48550/arXiv.1403.3741.
 - Shubham Pateria, Budhitama Subagdja, Ah-hwee Tan, and Chai Quek. Hierarchical reinforcement learning: A comprehensive survey. *ACM Computing Surveys (CSUR)*, 54(5):1–35, 2021.
 - Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Böhmer, and Shimon Whiteson. Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information Processing Systems*, 34:12208–12221, 2021.
 - Alexandre Piché, Valentin Thomas, Cyril Ibrahim, Yoshua Bengio, and Chris Pal. Probabilistic planning with sequential monte carlo methods. In *International Conference on Learning Representations*, 2018.
 - Thomas Pierrot, Valentin Macé, Jean-Baptiste Sevestre, Louis Monier, Alexandre Laterre, Nicolas Perrin, Karim Beguir, and Olivier Sigaud. Factored action spaces in deep reinforcement learning. 2021.
 - Yuta Saito, Jihan Yao, and Thorsten Joachims. Potec: Off-policy learning for large action spaces via two-stage policy decomposition. *arXiv* preprint arXiv:2402.06151, 2024.
 - Yunhao Tang and Shipra Agrawal. Discretizing continuous action space for on-policy optimization. In *Proceedings of the aaai conference on artificial intelligence*, volume 34, pp. 5981–5988, 2020.
 - Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
 - Arash Tavakoli, Fabio Pardo, and Petar Kormushev. Action branching architectures for deep reinforcement learning. 2017.

- Arash Tavakoli, Fabio Pardo, and Petar Kormushev. Action branching architectures for deep rein-forcement learning. In Proceedings of the aaai conference on artificial intelligence, volume 32, 2018. Tom Van de Wiele, David Warde-Farley, Andriy Mnih, and Volodymyr Mnih. Q-learning in enormous action spaces via amortized approximate maximization. arXiv preprint arXiv:2001.08116, 2020. Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International conference on machine learning*, pp. 3540–3549. PMLR, 2017. Zhihai Wang, Jie Wang, Qi Zhou, Bin Li, and Houqiang Li. Sample-efficient reinforcement learning
 - via conservative model-based actor-critic. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pp. 8612–8620, 2022.
 - Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control: Improved data-augmented reinforcement learning. arXiv preprint arXiv:2107.09645, 2021a.
 - Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images. In AAAI, pp. 10674–10681, 2021b.
 - Mingde Zhao, Zhen Liu, Sitao Luan, Shuyuan Zhang, Doina Precup, and Yoshua Bengio. A consciousness-inspired planning agent for model-based reinforcement learning. Advances in neural information processing systems, 34:1569–1581, 2021.

APPENDIX

The technical appendix includes:

- Sec. A: An overview of the action space sizes for all visual control tasks used in our experiments.
- Sec. B: Additional quantitative and qualitative results, including:
 - Sec. B.1: Visualization of policy evaluations.
 - Sec. B.2: Comparison of the planning efficiency per environment interaction step with TD-MPC2, which is based on MPPI planning.
 - Sec. B.3: Visualizations of the learned action-dimensional preferences, demonstrating that the high-level policy generates focused and interpretable control-unit prompts.
 - Sec. B.4: Additional results of FOCUS on low-dimensional tasks, highlighting its versatility across diverse visual control settings.
- Sec. C: A detailed description of the world model architecture and its associated learning objectives.
- Sec. D: The description of LLM usage in this work.

A SUMMARY OF EXPERIMENTAL ENVIRONMENTS

Table 1 summarizes the environmental setups in the DeepMind Control and MyoSuite benchmarks. While FOCUS is designed to tackle high-dimensional problems with large action spaces, its results on low-dimensional tasks highlight broad applicability.

Table 1: Overview of action space dimensions for all continuous control tasks in our experiments.

Environment	Tasks	Action Space Size (\mathbb{R}^d)	
High-dimensional DMC	Humanoid Stand	21	
	Humanoid Walk	21	
	Dog Walk	38	
	Dog Trot	38	
MyoSuite	Key Turn Hard	39	
	Pen Twirl Hard	39	
	Pose Hard	39	
Low-dimensional DMC	Walker Walk	6	
	Cheetah Run	6	
	Quadruped Walk	12	

B ADDITIONAL QUANTITATIVE AND QUALITATIVE RESULTS

B.1 VISUAL EVALUATION OF LEARNT POLICIES

We evaluate trained policies of different models on the DMC and Myosuite tasks and visualize episode frames for comparison. The full results on HumanoidWalk, DogTrot and KeyTurnHard are presented in Figure 7, Figure 8 and Figure 9 respectively.

B.2 PLANNING EFFICIENCY COMPARISON WITH TD-MPC2

We compare the planning efficiency of FOCUS and TD-MPC2 in Figure 10, measuring the time consumed from receiving an observation to producing an action. The experiments were conducted on a system equipped with an NVIDIA RTX 4090 GPU and a dual-socket Intel Xeon Gold 6240C CPU at 2.60GHz. For a fair comparison, both compared models are implemented in PyTorch 2.4.1, whereas the proposed FOCUS is implemented in JAX for the main text and all other experiments.

Although both methods rely on online planning, TD-MPC2 uses the *Model Predictive Path Integral* (MPPI) algorithm which is generally more computationally intensive than the proposed Sequential

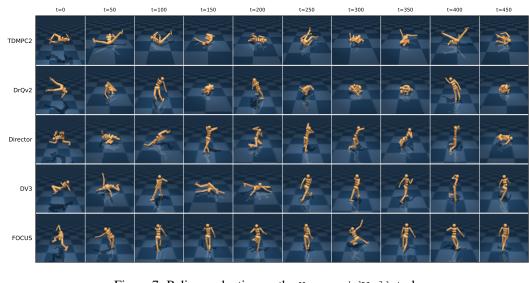


Figure 7: Policy evaluation on the HumanoidWalk task.

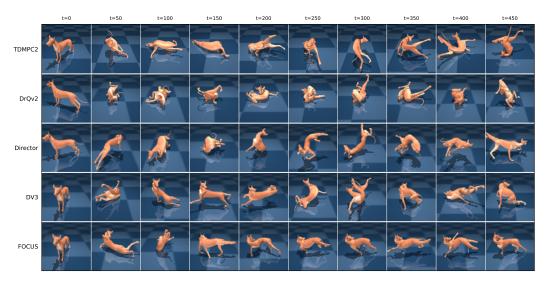


Figure 8: Policy evaluation on the DogTrot task.

Monte Carlo (SMC) planning employed in FOCUS. As shown, FOCUS achieves significantly lower latency than TD-MPC2, with a notably smaller increase in latency across varying planning horizons. This also highlights the scalability of FOCUS for long-horizon tasks.

B.3 VISUALIZATION OF LEARNED HIGH-LEVEL PREFERENCES

We visualize the preference prompts sampled by the learned high-level policy in Figure 11, and present the full visualization video in supplementary materials. The results show that our hierarchical policy can selectively prioritize action dimensions corresponding to key joints, such as the legs and arms in the humanoid tasks.

For instance, in the second row in Figure 11, the policy focuses on the *left knee* and the *x-actuator* at the right ankle to support the stand-up motion. A closer look at joint control axes reveals that the selected x-actuator at the right ankle is oriented vertically to the ground, making it particularly effective for initiating a standing-up behavior. In the third row, upon detecting that the right lower leg requires adjustment, the policy prioritizes control of the right knee joint to enable stable locomotion.

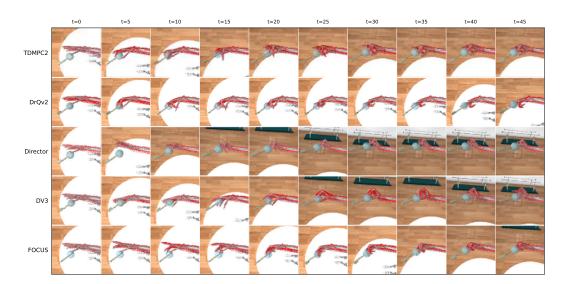


Figure 9: Policy evaluation on the KeyTurnHard task.

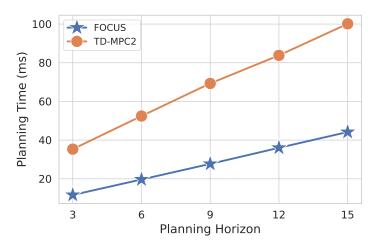


Figure 10: Time per interaction step to compare the efficiency of different planning algorithms.

B.4 APPLICABILITY TO LOW-DIMENSIONAL TASKS

A key concern regarding the proposed FOCUS is whether the high-level preference policy might hinder sufficient exploration in environments with small action spaces. To examine this, we evaluate FOCUS on four tasks from the DeepMind Control Suite characterized by low-dimensional action spaces: Walker Walk, Cheetah Run, and Quadruped Walk, while additional task details can be found in Table 1. As shown in Table 2, FOCUS achieves performance comparable to DreamerV3 across these low-dimensional control tasks, demonstrating its robustness and broad applicability.

C WORLD MODEL DETAILS

The DreamerV3 world model is implemented as a recurrent state-space model (RSSM), which maintains a deterministic hidden state h_t and a stochastic latent variable z_t to encode observations and predict future states. In our main text, we denote the concatenated feature $s_t = [h_t, z_t]$ as the

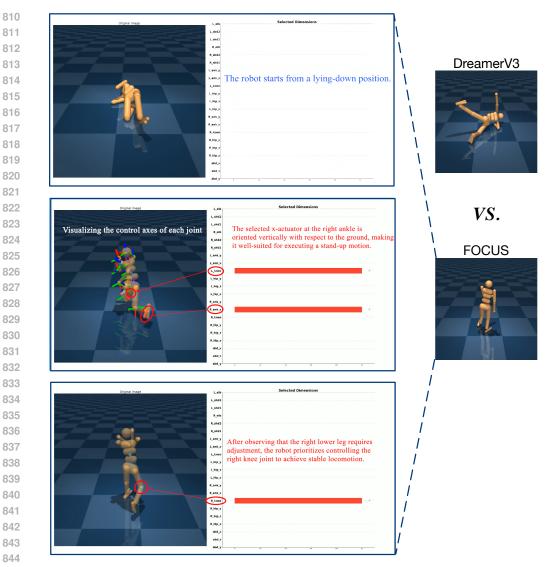


Figure 11: Visualization of learned action-dimensional preferences. Red bars indicate the action dimensions selected by the high-level policy. Please refer to our supplementary video for details.

Table 2: Episode return on low-dimensional DMC tasks. The comparable performance to DreamerV3 highlights the applicability of FOCUS across diverse continuous visual control setups.

Method	Walker Walk	Cheetah Run	Quadruped Walk	Average
FOCUS DreamerV3	973.3±13.5 966.4+21.19	910.1±61.44 901.9±26.2	953.0±8.33 943.4+29.1	945.5 937.2
TD-MPC2	332.9 ± 31.9	202.7 ± 2.3	162.2 ± 49.7	232.6

overall state representation. The world model can be formulated as follows:

Sequence model: $h_t = f_{\phi}(h_{t-1}, z_{t-1}, a_{t-1}),$ Encoder: $z_t \sim q_{\phi}(z_t \mid h_t, x_t)$, Dynamics predictor: $\hat{z}_t \sim p_{\phi}(\hat{z}_t \mid h_t)$, (7)Reward predictor: $\hat{r}_t = r_{\phi}(h_t, z_t)$, Continue predictor: $\hat{c}_t = c_{\phi}(h_t, z_t)$, Decoder: $\hat{x}_t = g_{\phi}(h_t, z_t),$

where f_{ϕ} is the recurrent transition function modeled by a GRU. q_{ϕ} and p_{ϕ} are implemented with MLPs. We use another MLP with separate heads r_{ϕ} and c_{ϕ} to predict the reward r_{t} and the episode continuation flag c_{t} , respectively. We use a CNN decoder to reconstruct the observation x_{t} .

The world model is trained by minimizing a variational objective that combines reconstruction losses with KL regularization. In particular, the total loss is written as:

$$\mathcal{L}(\phi) = \mathbb{E}_{\tau \sim \mathcal{B}} \left[\sum_{t=1}^{T} \left(\mathcal{L}_{t}^{\text{pred}} + \mathcal{L}_{t}^{\text{dyn}} + \mathcal{L}_{t}^{\text{rep}} \right) \right],$$

$$\mathcal{L}_{t}^{\text{pred}} = -\ln p_{\phi}(x_{t} \mid h_{t}, z_{t}) - \ln p_{\phi}(r_{t} \mid h_{t}, z_{t}) - \ln p_{\phi}(c_{t} \mid h_{t}, z_{t}),$$

$$\mathcal{L}_{t}^{\text{dyn}} = \text{KL} \left(\text{sg}(q_{\phi}(z_{t} \mid h_{t}, x_{t})) \parallel p_{\phi}(\hat{z}_{t} \mid h_{t}) \right),$$

$$\mathcal{L}_{t}^{\text{rep}} = \text{KL} \left(p_{\phi}(\hat{z}_{t} \mid h_{t}) \parallel \text{sg}(q_{\phi}(z_{t} \mid h_{t}, x_{t})) \right).$$
(8)

where $\mathcal{L}_t^{\text{pred}}$ is a self-supervised loss that reconstructs x_t and predicts the reward r_t and terminal signal c_t ; $\mathcal{L}_t^{\text{dyn}}$ is a KL divergence aligning the prior $p_{\phi}(\hat{z}_t \mid h_t)$ to the posterior $q_{\phi}(z_t \mid h_t, x_t)$; and $\mathcal{L}_t^{\text{rep}}$ is a reverse KL that regularizes the latent representations.

D USE OF LLMS

Large language models (LLMs) were used exclusively for language editing, including grammar correction and sentence polishing. They were not used for developing ideas, designing methods, conducting experiments, or interpreting results. Specifically, we provided prompts such as "Help polish this sentence to improve clarity and fluency in an academic writing style.", "Rewrite this sentence in a different expression", or "What is the alternative word for xxx". The output was then evaluated and further edited to stay close to this paper's writing style.