

V-Simba: Unleashing the Architectural Potential of RL in Visual Continuous Control

Anonymous authors

Paper under double-blind review

Keywords: Deep RL, Visual RL, Neural Network Design, Plasticity, Normalization

Summary

Improving sample efficiency remains a core challenge in reinforcement learning (RL), especially in real-world settings like robotics, where data collection is costly. This challenge is pronounced in visual RL, where high-dimensional inputs often obscure learning signals. While prior work in visual RL has focused on algorithmic solutions, such as better dynamics models or exploration strategies, recent advances in state-based RL show that architectural design alone can lead to significant gains in sample efficiency. This raises an important question: *Can these architectural principles transfer to visual RL?* In response, we introduce **V-Simba**, a simple yet effective visual RL architecture inspired by the Simba architecture from state-based RL. Built on top of Soft Actor-Critic (SAC) with data augmentation, V-Simba modifies the architecture by adding normalization layers to stabilize training and using pointwise convolutions to reduce computation. Despite its simplicity, V-Simba matches or outperforms the state-of-the-art methods across DMC, Adroit, and Meta-World benchmarks, while being more computationally efficient than DrQ-v2.

Contribution(s)

1. We identify and analyze severe training instabilities—such as sharp loss landscapes, dormant units, and feature collapse—that are inherent in standard convolutional architectures widely adopted in visual reinforcement learning.
Context: While previous visual RL research has primarily focused on algorithmic innovations, the underlying neural architectures have largely remained simple, often adopting the standard DrQ-v2 baseline. Our analysis reveals the vulnerabilities and capacity loss associated with these heavily utilized architectures.
2. We propose V-Simba, a novel visual RL architecture that integrates normalization layers, weight regularization, and a distributional critic to stabilize training, alongside large-stride and pointwise convolutions to maintain computational efficiency.
Context: Inspired by the recent Simba architectures in state-based RL that use principled architectural guidelines to constrain feature and parameter growth, we adapt these concepts to provide a suitable inductive bias specifically tailored for visual domains.
3. We experimentally demonstrate that V-Simba matches or outperforms state-of-the-art visual RL methods across 29 tasks spanning the DeepMind Control Suite, Adroit, and Meta-World benchmarks.
Context: We evaluate V-Simba against strong algorithmic baselines including DrQ-v2, MR.Q, TD-MPC2, DrM, and TACO, showing that robust architectural design alone can achieve superior sample and compute efficiency without introducing complex algorithmic add-ons.

V-Simba: Unleashing the Architectural Potential of RL in Visual Continuous Control

Anonymous authors

Paper under double-blind review

Abstract

1 Improving sample efficiency remains a core challenge in reinforcement learning (RL),
 2 especially in real-world settings like robotics, where data collection is costly. This chal-
 3 lenge is pronounced in visual RL, where high-dimensional inputs often obscure learning
 4 signals. While prior work in visual RL has focused on algorithmic solutions, such as bet-
 5 ter dynamics models or exploration strategies, recent advances in state-based RL show
 6 that architectural design alone can lead to significant gains in sample efficiency. This
 7 raises an important question: *Can these architectural principles transfer to visual RL?*
 8 In response, we introduce **V-Simba**, a simple yet effective visual RL architecture in-
 9 spired by the Simba architecture from state-based RL. Built on top of Soft Actor-Critic
 10 (SAC) with data augmentation, V-Simba modifies the architecture by adding normaliza-
 11 tion layers to stabilize training and using pointwise convolutions to reduce computation.
 12 Despite its simplicity, V-Simba matches or outperforms the state-of-the-art methods
 13 across DMC, Adroit, and Meta-World benchmarks, while being more computationally
 14 efficient than DrQ-v2.

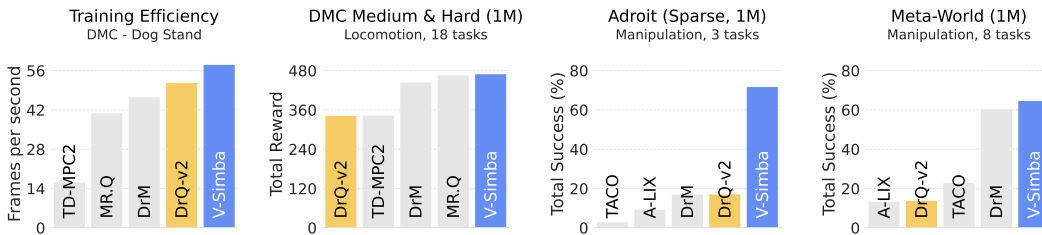


Figure 1: **Benchmark Summary.** We evaluate the effectiveness of V-Simba across 29 visual continuous control tasks spanning multiple domains, with a *single* set of hyperparameters. By incorporating V-Simba into Soft Actor-Critic with data augmentation, it matches or outperforms state-of-the-art visual RL methods, demonstrating better sample and compute efficiency.

15 1 Introduction

16 Deep reinforcement learning (RL) has long been a prominent approach for solving continuous control tasks. However, RL typically relies on an extensive amount of trial-and-error within the environment, which is often expensive in terms of time, compute, and real-world constraints. This issue is exacerbated in visual RL, where agents must learn from high-dimensional, noisy, and often partially observable image inputs. Consequently, improving sample efficiency (i.e., learning effectively from limited interaction data) has become a key research topic in visual RL.

22 To improve sample efficiency, recent work has largely concentrated on algorithmic innovations, including enhanced representation learning (Yarats et al., 2021b), latent dynamics modeling (Fujimoto et al., 2025; Zheng et al., 2023), world models (Hansen et al., 2023; Hafner et al., 2023),

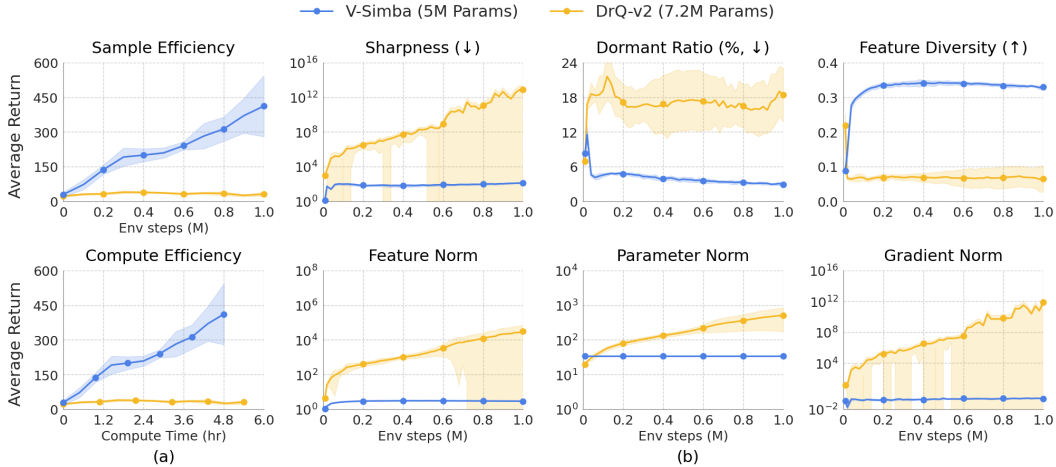


Figure 2: **DrQ-v2 v.s. V-Simba**. Comparison of the DrQ-v2 architecture and our V-Simba in the `Dog Stand` environment. Both are evaluated on Soft Actor-Critic (SAC) with results averaged over 5 seeds. **(a)** V-Simba is substantially more efficient than DrQ-v2 in terms of both sample and compute. **(b)** Unlike DrQ-v2, V-Simba has stable learning dynamics, indicated by smooth loss landscape, low dormant ratio, high feature diversity, and well-controlled feature, parameter, and gradient norms. Detailed explanations of each metric are provided in Appendix 7.2.

25 and improved exploration strategies (Burda et al., 2018b; Xu et al., 2023). Yet, despite these ad-
 26 vances, the underlying neural architectures have remained relatively simple. A prominent example
 27 is DrQ-v2 (Yarats et al., 2021a), which combines the DDPG algorithm (Lillicrap, 2015) with data
 28 augmentation (Laskin et al., 2020). Its architecture consists of a shallow convolutional encoder, fol-
 29 lowed by a large fully connected layer and a single layer normalization layer (Lei Ba et al., 2016)
 30 in-between. Due to its simplicity and strong empirical performance, DrQ-v2 has become the de facto
 31 standard in visual RL, and many state-of-the-art methods (Xu et al., 2023; Zheng et al., 2023; Cetin
 32 et al., 2022; Sukhija et al., 2024) adopt DrQ-v2’s architecture with minimal modifications.

33 However, our analysis reveals that this commonly adopted architecture suffers from severe training
 34 instabilities. As shown in the top row of Figure 2.(b), DrQ-v2 exhibits sharp loss landscapes that
 35 correlate with poor generalization (Foret et al., 2020; Lee et al., 2024a), a high fraction of dormant
 36 units representing plasticity loss (Sokar et al., 2023a), and low feature diversity indicating feature
 37 collapse (Woo et al., 2023). In contrast, recent advances in state-based RL (Lee et al., 2024b; Bhatt
 38 et al., 2024; Lee et al., 2025; Palenicek et al., 2025) demonstrate that carefully designed architectures
 39 can effectively mitigate these instabilities. Notably, the Simba series of architectures (Lee et al.,
 40 2024b; 2025) introduce principled architectural guidelines that stabilize training by constraining the
 41 growth of features, weights, and gradients through targeted normalization and regularization.

42 While effective in state-based tasks, the Simba architecture lacks a suitable inductive bias for visual
 43 data. In response, we propose **V-Simba**, a simple yet effective architecture for visual continuous
 44 control, as a means of applying the underlying principles from the Simba model series to visual RL
 45 domains. Built on top of Soft Actor-Critic (SAC) (Haarnoja et al., 2018), V-Simba incorporates three
 46 core architectural components: (1) normalization layers (LN) to control feature norms, (2) ℓ_2 weight
 47 regularization to limit parameter growth, and (3) a distributional critic with reward normalization
 48 to stabilize gradients. To ensure computational efficiency, V-Simba applies early downsampling via
 49 large-stride convolutions and makes extensive use of lightweight pointwise convolutions (Hua et al.,
 50 2018) and max-pooling operations (Krizhevsky et al., 2012).

51 We evaluate V-Simba on three standard benchmarks: DMControl (Tassa et al., 2018), Adroit (Ra-
 52 jeswaran et al., 2017), and Metaworld (Yu et al., 2020) using a single set of hyperparameters across
 53 all tasks. Despite its simplicity, V-Simba consistently outperforms DrQ-v2, while reducing both

54 model size (7.2M \rightarrow 5.0M) and training time (5.4 \rightarrow 4.8 hours for 1M DMControl steps). Moreover,
 55 V-Simba is competitive with leading vision-based methods, matching or surpassing MR.Q (Fujimoto
 56 et al., 2025) and TD-MPC2 (Hansen et al., 2023) in DMControl and outperforming DrM (Xu et al.,
 57 2023) and TACO (Zheng et al., 2023) in dexterous manipulation tasks.

58 V-Simba is intended to offer a strong, stable, and efficient architectural foundation for advancing
 59 visual continuous control. We hope our work highlights the untapped potential of principled archi-
 60 tecture design within the visual RL community.

61 2 Related Work

62 Learning solely from high-dimensional visual observations poses significant challenges in RL. Due
 63 to the partially observable nature of the observation space (Section 3.1), visual RL agents suffer from
 64 poor sample efficiency and large generalization gaps compared to their state-based counterparts (Ma
 65 et al., 2022).

66 **Algorithmic approaches for visual RL** have primarily focused on: (1) representation learning
 67 via auxiliary tasks predicting future latent states, either as auxiliary losses in model-free meth-
 68 ods (Stooke et al., 2021; Zheng et al., 2023; Schwarzer et al., 2020; 2021; Kim et al., 2022; Lee et al.,
 69 2020; Van Hoof et al., 2016; Yu et al., 2022; Gelada et al., 2019; Seo et al., 2022; Ni et al., 2024; Yu
 70 et al., 2021; Fujimoto et al., 2021; McInroe et al., 2021; Fujimoto et al., 2025) or separate dynam-
 71 ics models in model-based methods (Hansen et al., 2023; Hafner et al., 2023; Lin et al., 2025; Wu
 72 et al., 2023; Ha & Schmidhuber, 2018; Finn et al., 2016; Watter et al., 2015); (2) data augmentation,
 73 especially random shifts (Laskin et al., 2020; Kostrikov et al., 2020; Yarats et al., 2021a), for better
 74 efficiency and generalization; and (3) exploration methods, including planning (Sekar et al., 2020;
 75 Wang et al., 2023), curiosity-driven (Pathak et al., 2017; Burda et al., 2018a; Guo et al., 2022), and
 76 information maximization (Sukhija et al., 2024). These algorithmic innovations have driven rapid
 77 progress in visual RL, leading to continuous improvements in sample efficiency.

78 **Architectural design for visual RL** has received comparatively little attention compared to algo-
 79 rithmic innovations. The field has largely maintained shallow convolutional neural network (CNN)
 80 architectures similar to the one established by DQN (Mnih et al., 2015) over a decade ago. While
 81 some works have incorporated architectural elements from computer vision—such as ResNet-like
 82 architectures in Impala (Espeholt et al., 2018), BBF (Schwarzer et al., 2023), and EfficientZero (Ye
 83 et al., 2021), or transformers in DTQN (Esslinger et al., 2022)—these modifications were often in-
 84 troduced alongside complex algorithmic methods. This entanglement has obscured the true contribu-
 85 tion of architectural design to performance improvements. Recent studies have identified the benefits
 86 of normalization techniques (Lyle et al., 2023; Ball et al., 2023; Lyle et al., 2024), but these have
 87 generally been applied to conventional CNN encoders with minimal architectural modifications. To
 88 the best of our knowledge, aside from a few isolated attempts such as adding global average pooling
 89 (Trumpp et al., 2025) or using Mixture-of-Experts (Obando-Ceron et al., 2024; Sokar et al., 2024),
 90 no substantial architectural innovations have been sufficiently explored in visual RL (Espeholt et al.,
 91 2018; Huang et al., 2025).

92 3 Preliminary

93 As a preliminary, we briefly describe the problem setup of visual RL, DrQ-v2 architecture (Yarats
 94 et al., 2021a), and the Soft Actor-Critic algorithm (Haarnoja et al., 2018), as these form the founda-
 95 tion for our proposed architecture.

96 3.1 Visual Reinforcement Learning

97 Reinforcement learning (RL) is typically formulated as a Markov Decision Process (MDP) (Bell-
 98 man, 1957), defined by the tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ of state space \mathcal{S} , action space \mathcal{A} , transition function
 99 $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$, reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and discount factor $\gamma \in [0, 1)$. From an

100 initial state $s_0 \in \mathcal{S}$, the objective is to find an optimal policy $\pi^* : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ that maximizes the
 101 expected discounted return $\mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$. Visual RL is a subclass of this problem where
 102 the agent does not have access to the true state $s \in \mathcal{S}$, but instead receives high-dimensional pixel
 103 observations $o \in \mathcal{O}$ of the system. Since these observations may not fully capture the true state, the
 104 problem is modeled as a Partially Observable Markov Decision Process (POMDP) (Bellman, 1957)
 105 represented by the tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, P, R, \gamma)$ where \mathcal{O} denotes the observation space.

106 3.2 Data-regularized Q-learning

107 Data-regularized Q-learning (DrQ-v2) (Yarats et al., 2021a) is a model-free RL algorithm that has
 108 emerged as a strong baseline in visual RL due to its simplicity, efficiency, and competitive perfor-
 109 mance. It builds upon the Deep Deterministic Policy Gradient (DDPG) algorithm (Lillicrap, 2015)
 110 by incorporating two key modifications: (1) extensive use of data augmentation via random shift
 111 transformations, and (2) target Q-function stabilization through exponential moving average (EMA)
 112 updates.

113 At its core, DrQ-v2 improves sample efficiency in off-policy learning by generating augmented
 114 views of each observation, thereby increasing data diversity. This augmentation acts as a regular-
 115 izer, mitigating overfitting to specific visual patterns. Despite its empirical effectiveness, DrQ-v2
 116 employs a notably lightweight architecture: a shallow convolutional encoder followed by an MLP-
 117 based prediction head, with a single normalization layer (Lei Ba et al., 2016) in between.

118 While DrQ-v2 has become the de facto architecture for many recent visual RL algorithms (Xu et al.,
 119 2023; Zheng et al., 2023), its architectural simplicity leaves room for improvement in stability and
 120 representational capacity.

121 3.3 Soft Actor-Critic (SAC)

122 Soft Actor-Critic (SAC) is a prominent off-policy algorithm for continuous control. It aims to max-
 123 imize both expected cumulative reward and policy entropy, where $\tau = (o, a, r, o')$ represents a
 124 transition tuple. SAC comprises a stochastic policy $\pi_\theta(a|o)$, a Q-function $Q_\phi(o, a)$, and an entropy
 125 coefficient α that balances reward maximization and entropy regularization. The policy network is
 126 optimized to maximize the expected return while encouraging exploration through entropy. This
 127 objective is formalized as:

$$\mathcal{L}_\pi = \mathbb{E}_{\bar{a} \sim \pi_\theta} [\alpha \log \pi_\theta(\bar{a}|o) - Q_\phi(o, \bar{a})]. \quad (1)$$

128 The Q-function $Q_\phi(o, a)$ is trained to minimize the Bellman residual:

$$\mathcal{L}_Q = (Q_\phi(o, a) - (r + \gamma Q_{\bar{\phi}}(o', a') - \alpha \log \pi_\theta(a'|o')))^2, \quad (2)$$

129 where $a' \sim \pi_\theta(\cdot|o')$, and $Q_{\bar{\phi}}$ represents the target Q-network updated via an exponential moving
 130 average of ϕ .

131 4 Method

132 V-Simba leverages architectural design from state-based RL to stabilize optimization dynamics and
 133 improve computational efficiency in visual RL. Our design follows two core principles: (1) stabiliz-
 134 ing optimization (Section 4.1), and (2) maintaining computational efficiency (Section 4.2). The final
 135 architecture builds on these principles (Section 4.3).

136 4.1 Design Philosophy I: Stabilizing Optimization

137 As shown in Figure 2, DrQ-v2 suffers from unstable optimization during training. While layer nor-
 138 malization (LN) (Lei Ba et al., 2016) and residual connections (He et al., 2020) effectively stabilize

139 supervised learning, visual RL methods often underuse them—DrQ-v2, for instance, employs only a
 140 single normalization layer without residuals. We incorporate both components to improve stability.

141 However when adding LayerNorm, one must consider its relationship with the gradient. Concretely,
 142 LayerNorm introduces scale invariance: for any scalar $c > 0$ and weight matrix W ,

$$\text{Norm}(cWx) = \text{Norm}(Wx), \tag{3}$$

143 which causes gradients to scale inversely with parameter magnitude:

$$\nabla_W \text{Norm}(cWx) = \frac{1}{c} \nabla_W \text{Norm}(Wx). \tag{4}$$

144 As parameter norms grow during training, gradients diminish, reducing learning ability (Lyle et al.,
 145 2024; Palenicek et al., 2025). Moreover, uneven growth across layers causes inconsistent gradient
 146 scales, destabilizing optimization (Lee et al., 2025). This highlights the importance of controlling
 147 weight and gradient norms, in addition to the feature norm. Thus, we employ the following design
 148 choices to achieve stable norms.

149 We first opt LayerNorm as the forefront layer of both encoder and critic module, in order to control
 150 the norm of not only their intermediate features but also their inputs. While unusual for convolu-
 151 tional networks, this resembles the Dual PatchNorm design (Kumar et al., 2023) which has been
 152 empirically shown to stabilize the gradients of embedding layer¹. For preventing parameter growth,
 153 we surprisingly found a simple ℓ_2 weight regularization to be sufficient, as shown in Figure 2.

154 We further stabilize gradients by employing a distributional critic with KL divergence loss (Belle-
 155 mare et al., 2017) and reward normalization (Lee et al., 2025). The KL divergence loss is more robust
 156 to noisy targets than mean squared error due to its smoother loss landscape (Farebrother et al., 2024),
 157 while reward normalization ensures consistent learning signals despite varying reward scales.

158 Specifically, reward normalization maintains unit variance in expected returns. Given reward r_t at
 159 time t , we track the discounted return:

$$G_t \leftarrow \gamma G_{t-1} + r_t \tag{5}$$

160 with G_t re-initialized to 0 at the start of each episode. Let $\sigma_{t,G}^2$ denotes the running variance of G_t .
 161 Each reward is then scaled as:

$$\bar{r}_t \leftarrow \frac{r_t}{\sqrt{\sigma_{t,G}^2 + \epsilon}}, \tag{6}$$

162 **4.2 Design Philosophy II: Maintaining Computational Efficiency**

163 Adding normalization layers and regularizations increases training cost, so reducing computation
 164 is crucial. We find that most of DrQ-v2’s computation cost comes from early convolutional lay-
 165 ers processing high-resolution inputs. We apply early downsampling via large-stride convolutions,
 166 a common practice in ResNet (He et al., 2016), ConvNeXt (Liu et al., 2022), and Vision Trans-
 167 former (Dosovitskiy et al., 2020). This results in an early reduction in spatial resolution and in turn,
 168 the computational cost of subsequent convolution layers.

169 We further cut computation by replacing convolutional layers into more cost-effective alternatives.
 170 For spatial convolutions, we instead utilize pointwise (1×1 kernel) convolutions (Hua et al., 2018),
 171 which operate channel-wise without mixing spatial information, preserving spatial details at a lower
 172 cost. When downsampling, instead of 2×2 kernel convolutional layers, which induces significantly
 173 less overhead in both training and inference. Empirically, we found that these alterations does not
 174 meaningfully alter the learning dynamics and the learning curves.

¹In practice, we adopt the shift-and-norm strategy introduced in SimbaV2 (Lee et al., 2025) to preserve magnitude infor-
 mation. We use ℓ_2 -norm for action inputs however, as when $|\mathcal{A}| = 1$ shift-and-LN always outputs $[-1, 1]$.

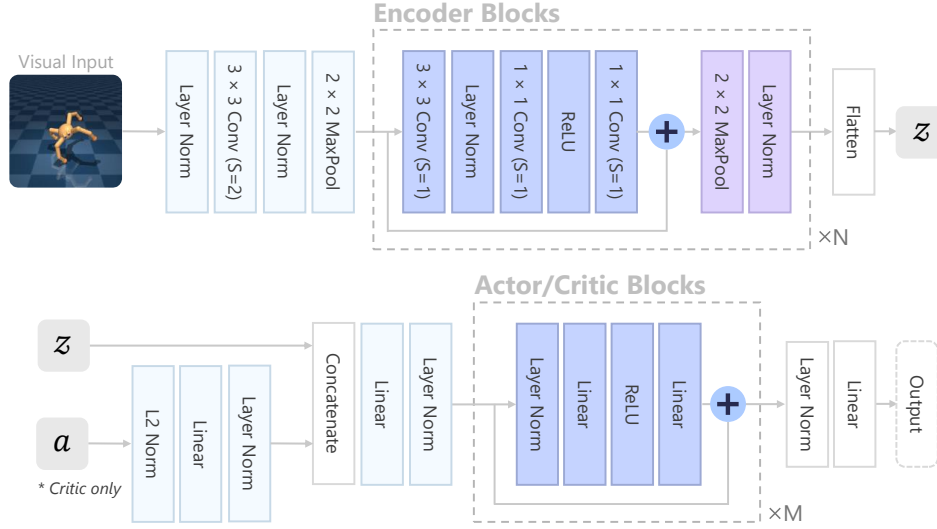


Figure 3: **V-Simba architecture.** We aim to develop an architecture that constrains its feature norm, weight norm, and gradient norm for better stability and generalization. Precisely, we make extensive use of layer normalization and residual connections for stable features and gradients, and incorporate point-wise convolutions in tandem with spatial convolutions for computational efficiency.

175 4.3 V-Simba Architecture

176 Building on our design principles, we now detail the V-Simba architecture (Figure 3).

177 **Image Preprocessing.** The input $o \in \mathbb{R}^{84 \times 84 \times 9}$ is a stack of the last three RGB frames. We first
 178 apply an alternating sequence of normalization and reduction steps. Specifically, the input passes
 179 through an initial LayerNorm, followed by a 3×3 convolution (stride 2), a second LayerNorm, and
 180 finally a 2×2 max-pooling layer (stride 2). This results in downsampled and normalized features
 181 $f_0 \in \mathbb{R}^{21 \times 21 \times 32}$.

182 **Encoder.** The encoder consists of two sequential blocks transforming and downsampling features:
 183 $f_0 \xrightarrow{\text{Block}_1} f_1 \xrightarrow{\text{Block}_2} f_2$, where $f_1 \in \mathbb{R}^{10 \times 10 \times 32}$, $f_2 \in \mathbb{R}^{5 \times 5 \times 32}$.

184 Each encoder block processes input f_i as follows:

- 185 1. A 3×3 convolution to aggregate spatial features without changing resolution.
- 186 2. Feature normalization (LayerNorm) to provide stable features to the subsequent layers.
- 187 3. Two pointwise (1×1) convolutions with nonlinearities to refine and filter features, followed by
 188 a residual connection for stable gradient flow. Here, we employ an inverted bottleneck with $4 \times$
 189 expansion, following ConvNext (Liu et al., 2022) and Simba (Lee et al., 2024b),
- 190 4. Downsampling with a 2×2 maxpool layer with stride 2 to reduce spatial resolution
- 191 5. Applying LayerNorm to normalize the features before passing to the next block

192 After the second block, f_2 is flattened into the latent state vector $z \in \mathbb{R}^{800}$.

193 **Predictor.** The latent vector z feeds into separate actor and critic heads, each followed by a linear
 194 layer and LayerNorm. For the critic, actions are separately embedded and concatenated with image
 195 embedding. The actor and critic embeddings have dimensions $z_\pi \in \mathbb{R}^{128}$, $z_Q \in \mathbb{R}^{512}$ following Lee
 196 et al. (2024b).

197 Each embedding passes through residual nonlinear blocks: one block for the actor and two for the
 198 critic. Finally, the actor output passes through LayerNorm, linear layer and tanh activation, while
 199 the critic output passes through LayerNorm and linear layer modeling the Q-value distribution.

200 **5 Experiments**

201 We now provide an empirical evaluation of V-Simba:

- 202 1. **Performance Evaluation** (Sections 5.2), comparing V-Simba against leading visual RL meth-
203 ods to demonstrate its effectiveness across diverse benchmarks.
- 204 2. **Ablation Study** (Section 5.3), conducting experiments demonstrating the contribution of each
205 architectural component in V-Simba.

206 **5.1 Experimental Setup**

207 **Environment.** We consider a total of 29 continuous control tasks spanning 3 benchmarks: Deep-
208 Mind Control (DMC) Suite (Tassa et al., 2018), Adroit (Rajeswaran et al., 2017), and Meta-
209 World (Yu et al., 2020). Figure 4 shows the visualization of each task. These environments pose
210 diverse challenges, including high-dimensional action spaces, sparse rewards, and complex dexter-
211 ous manipulation, often under rich visual observations with shading and textures. Consequently, to
212 solve the tasks, prior visual RL methods typically require either large volumes of frames or privi-
213 leged information such as low-level robot states.

214 **Baselines.** In experiments, we compare V-Simba against a diverse set of state-of-the-art visual
215 RL methods exemplifying three key algorithmic strategies: data and model regularization (DrQ-
216 v2 (Yarats et al., 2021a), A-LIX (Cetin et al., 2022)), advanced exploration (DrM (Xu et al., 2023)),
217 and model-based representation learning (TACO (Zheng et al., 2023), TD-MPC2 (Hansen et al.,
218 2023), MR.Q (Fujimoto et al., 2025)). Notably, A-LIX, TACO, and DrM build upon DrQ-v2 (see
219 Section 3.2): A-LIX stabilizes training by adaptively regularizing the encoder’s gradients; TACO
220 leverages a latent dynamics loss for richer representations; and DrM integrates dormant ratio (Sokar
221 et al., 2023b)-guided mechanisms that balance exploration-exploitation dynamically. While these
222 variants benefit from task-specific hyperparameter tuning, our method uses the *same* hyperparam-
223 eters across all tasks. Whenever possible, we report original paper results; otherwise, we run the
224 authors’ official implementations.

225 **5.2 Performance Evaluation**

226 **DMC Medium.** We begin by evaluating V-Simba on DMC Medium, consisting of 11 mid-difficulty
227 tasks from DMC. As shown in Figure 6, our base algorithm, DrQ-v2, falls behind model-based
228 methods such as TD-MPC2 and MR.Q. However, simply replacing DrQ-v2’s neural network with
229 our proposed architecture, V-Simba, yields substantial performance gains. As a result, V-Simba
230 surpasses TD-MPC2 and achieves results competitive with leading algorithm, MR.Q, highlighting
231 the impact of architectural improvements.

232 **DMC Hard.** We further assess V-Simba on DMC Hard, a set of 7 high-difficulty tasks in DMC,
233 characterized by complex kinematics and high-dimensional control. Figure 5 shows that V-Simba
234 performs competitively with MR.Q, though full task success remains elusive. We believe that this

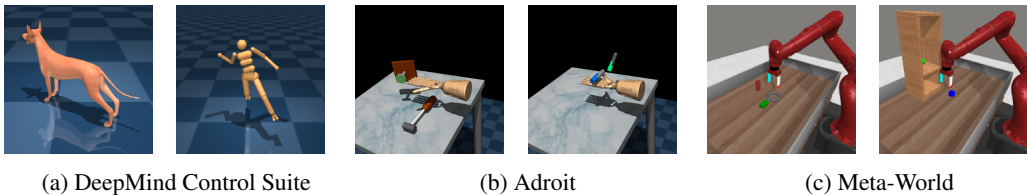


Figure 4: **Environment Visualization.** We evaluate our V-Simba on 3 visual continuous control benchmarks: DeepMind Control Suite (Tassa et al., 2018), Adroit (Rajeswaran et al., 2017), and Meta-World (Yu et al., 2020).

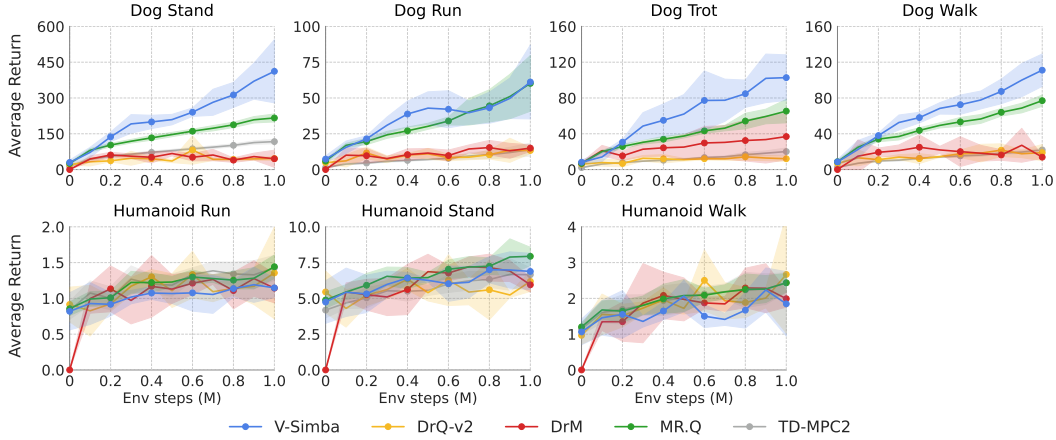


Figure 5: **DeepMind Control Suite - Hard.** Average episode returns on 7 hard-level tasks from DeepMind Control Suite (Tassa et al., 2018). Each curve represents the mean performance across 5 random seeds per algorithm; shaded areas indicate 95% bootstrap confidence intervals.

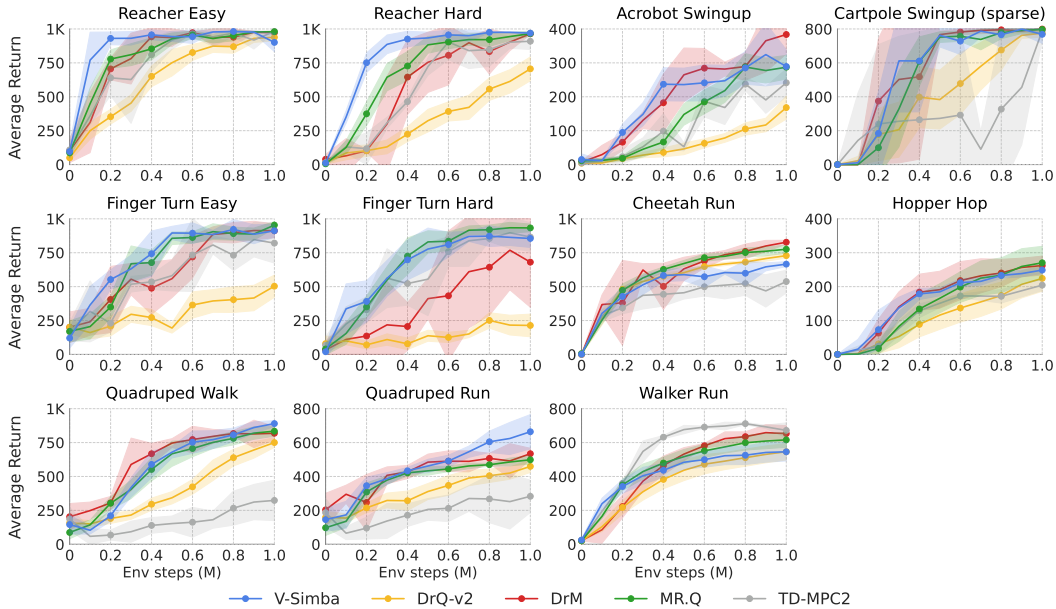


Figure 6: **DeepMind Control Suite - Medium.** Average episode returns on 11 medium-level tasks from DeepMind Control Suite (Tassa et al., 2018). Each curve represents the mean performance across 5 random seeds per algorithm; shaded areas indicate 95% bootstrap confidence intervals.

235 observation suggests that concurrent advances in both algorithm design and architectural representa-
 236 tion are needed in visual RL, to close the gap with state-based performance.

237 **Adroit - Sparse.** Moving to more intricate scenarios, we evaluate V-Simba on Adroit under the
 238 challenging *sparse-reward* setting. In this domain, the agent must control a dexterous hand-arm
 239 system to perform complex manipulation such as opening a door or using tools like a hammer.
 240 These tasks pose significant challenges for visual RL, often requiring over 5 million environment
 241 frames and access to privileged robot state inputs for successful learning. The comparison results
 242 are shown in Figure 7. V-Simba reliably solves or approaches solving all tasks using only 1 million
 243 frames. In contrast, DrM—the previous state-of-the-art method—fails to learn meaningful behavior,
 244 despite leveraging privileged state vectors. Notably, V-Simba is the only method to solve Hammer

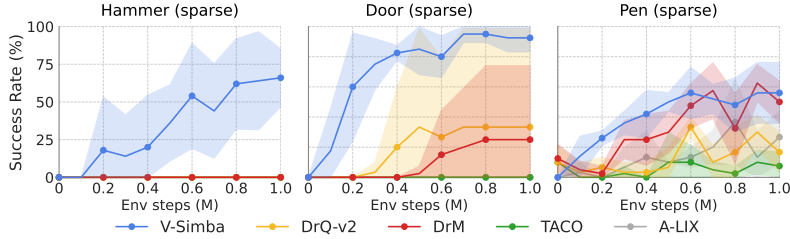


Figure 7: **Adroit - Sparse.** Average success rates on 3 sparse-reward tasks from Adroit (Rajeswaran et al., 2017). Each curve represents the mean performance across 5 random seeds per algorithm; shaded areas indicate 95% bootstrap confidence intervals.

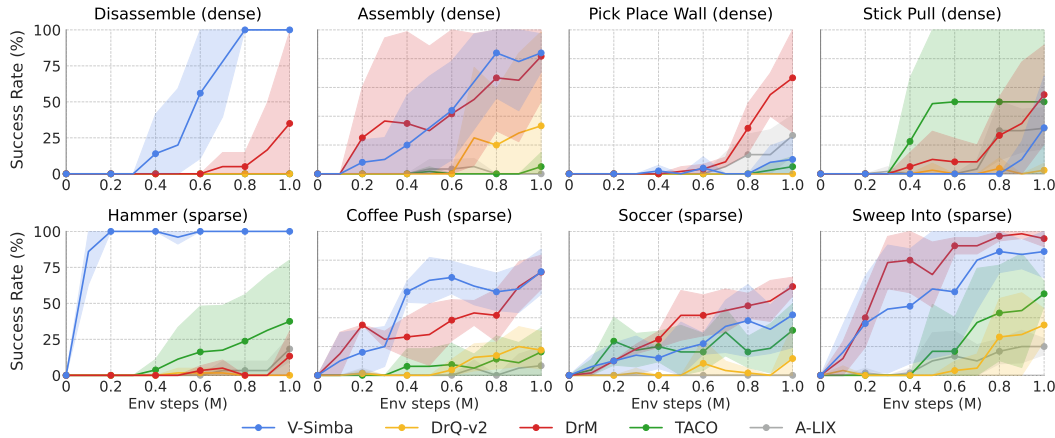


Figure 8: **Meta-World.** Average success rates on 8 tasks from the Meta-World (Yu et al., 2020). Each curve represents the mean performance across 5 random seeds per algorithm; shaded areas indicate 95% bootstrap confidence intervals.

245 with 1 million environment steps. These results underscore V-Simba’s strong sample-efficiency and
 246 effectiveness in high-dimensional visual control settings.

247 **Meta-World.** We also benchmark V-Simba on Meta-World, which demands precise object ma-
 248 nipulation. We consider 4 medium-difficulty tasks: *Coffee Push*, *Soccer*, *Sweep Into*,
 249 and *Hammer*, and 4 high-difficulty tasks: *Assembly*, *Stick Pull*, *Pick Place Wall*, and
 250 *Disassemble*. For the medium tasks, we adopt a sparse-reward setting by replacing the ground-
 251 truth reward functions with binary success signals, following Yu et al. (2020), to increase task dif-
 252 ficulty. As shown in Figure 8, while DrQ-v2 struggles to learn in most tasks, V-Simba significantly
 253 improves over DrQ-v2 and matches or surpasses leading baselines, demonstrating superior sample
 254 efficiency. A notable performance improvement can be seen in the *Disassemble* and *Hammer*
 255 task, where V-Simba was able to consistently achieve almost perfect success rate, whereas prior
 256 works have failed in few trials.

257 **5.3 Ablation Study**

258 To assess the impact of each component on V-Simba’s overall performance, we evaluate variants
 259 that remove or modify one component at a time. The results are reported in Table 1.

260 We first investigate the effect of normalization layers (Table 1.(a)-(c)). **No Normalization Layers**
 261 removes LayerNorm entirely from the network, whereas **No Input Normalization** only removes two
 262 LayerNorms: for image and action inputs in encoder and critic respectively. **LN w/o γ, β** removes

Table 1: **Ablation Study.** We exclude or modify each component in V-Simba and report their final performance on each benchmark, averaged over 3 random seeds. Each cell is highlighted base on their relative percentile difference to V-Simba, namely: positive (> 0.01), mildly negative $[-0.05, -0.01]$, damaging $[-0.1, -0.05]$, and catastrophic $[-1.0, -0.1]$.

Ablation	DMC (18) Return (1k)	Adroit (3) Success Rate	MetaWorld (8) Success Rate	All (29) -
Normalization Layers				
(a) No Normalization Layers	0.442 \pm 0.095	0.694 \pm 0.122	0.497 \pm 0.167	0.492 \pm 0.081
(b) No Input Normalization	0.401 \pm 0.098	0.612 \pm 0.133	0.642 \pm 0.158	0.502 \pm 0.083
(c) LN \rightarrow LN w/o γ, β	0.453 \pm 0.073	0.668 \pm 0.133	0.586 \pm 0.125	0.522 \pm 0.063
Residual and Weight Decay				
(d) No Residual Connection	0.453 \pm 0.098	0.623 \pm 0.178	0.618 \pm 0.154	0.526 \pm 0.081
(e) No Weight Decay	0.452 \pm 0.099	0.591 \pm 0.144	0.513 \pm 0.171	0.492 \pm 0.080
Value Learning				
(f) No Categorical Critic	0.441 \pm 0.096	0.591 \pm 0.144	0.576 \pm 0.167	0.504 \pm 0.079
(g) No Reward Scaling	0.439 \pm 0.094	0.683 \pm 0.144	0.599 \pm 0.158	0.519 \pm 0.080
V-Simba	0.467 \pm 0.075	0.725 \pm 0.167	0.668 \pm 0.115	0.549 \pm 0.060

263 the bias and scale parameters of LayerNorm. In summary, by removing certain normalization layers
264 or components, the network loses the control over the feature norms, leading to degradation.

265 Next, we quantify the importance of residual connections and weight decay (Table 1.(d)-(e)). Both
266 residual connection and weight decay, along with their well-known benefits, are also known to bias
267 the network towards simple solutions for improved robustness (Teney et al., 2024; Lee et al., 2024b).
268 Removing such components led to visible drop in performance, similar to removing normalization
269 layers.

270 Finally, categorical critic and reward scaling are critical components, as they reformulate the re-
271 gression problem into a categorical prediction, giving a much more stable gradient and learning
272 dynamics. Reverting back to regression loss led to diminished performance (Table 1.(f)). Even with
273 categorical loss, leaving no bounds to the reward scales led to similar consequences (Table 1.(g)),
274 highlighting the importance of assuring the Q-values to stay in a certain range.

275 6 Lessons and Opportunities

276 In this work, we introduce V-Simba, a simple yet effective neural network architecture for visual
277 continuous control, inspired by the Simba architecture from state-based RL (Lee et al., 2024b). By
278 combining feature normalization, weight regularization, and a distributional critic, V-Simba achieves
279 superior performance over prior visual RL methods across multiple benchmarks with minimal al-
280 gorithmic changes. Additionally, V-Simba reduces computational cost by integrating early down-
281 sampling through large-stride convolution and pointwise convolution layers, enabling faster training
282 than DrQ-v2 (Yarats et al., 2021a). We believe our work does not oppose the current trend of adopt-
283 ing model-based learning or exploration strategies; rather, it offers a complementary approach that
284 can be integrated with subsequent studies.

285 Moreover, in recent years, reinforcement learning for robotic control has gained increased atten-
286 tion. However, limited sample efficiency remains a significant barrier to real-world adoption. While
287 simulators provide valuable virtual environments (Makoviychuk et al., 2021; Zakka et al., 2025),
288 rendering high-resolution images with complex object interactions is still computationally expen-
289 sive and difficult to parallelize. This underscores the importance of improving sample efficiency.
290 V-Simba offers a lightweight architectural solution using well-established components that are easy
291 to integrate into existing algorithms. Its simplicity allows practitioners to adopt and extend it with
292 minimal overhead. We hope V-Simba serves as an architectural foundation to accelerate progress in
293 the robotics community.

294 **References**

- 295 Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learn-
296 ing with offline data. In *International Conference on Machine Learning*, pp. 1577–1594. PMLR,
297 2023.
- 298 Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement
299 learning. In *International conference on machine learning*, pp. 449–458. PMLR, 2017.
- 300 Richard Bellman. A markovian decision process. *Journal of mathematics and mechanics*, pp. 679–
301 684, 1957.
- 302 Aditya Bhatt, Daniel Palenicek, Boris Belousov, Max Argus, Artemij Amiranashvili, Thomas Brox,
303 and Jan Peters. Crossq: Batch normalization in deep reinforcement learning for greater sample
304 efficiency and simplicity. In *The Twelfth International Conference on Learning Representations*,
305 2024. URL <https://openreview.net/forum?id=PczQtTsTIX>.
- 306 Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
307 distillation. *arXiv preprint arXiv:1810.12894*, 2018a.
- 308 Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
309 distillation. *arXiv preprint arXiv:1810.12894*, 2018b.
- 310 Edoardo Cetin, Philip J Ball, Steve Roberts, and Oya Celiktutan. Stabilizing off-policy deep rein-
311 forcement learning from pixels. *arXiv preprint arXiv:2207.00986*, 2022.
- 312 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
313 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
314 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
315 *arXiv:2010.11929*, 2020.
- 316 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron,
317 Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance
318 weighted actor-learner architectures. In *International conference on machine learning*, pp. 1407–
319 1416. PMLR, 2018.
- 320 Kevin Esslinger, Robert Platt, and Christopher Amato. Deep transformer q-networks for partially
321 observable reinforcement learning. *arXiv preprint arXiv:2206.01078*, 2022.
- 322 Jesse Farebrother, Jordi Orbay, Quan Vuong, Adrien Ali Taïga, Yevgen Chebotar, Ted Xiao, Alex
323 Irpan, Sergey Levine, Pablo Samuel Castro, Aleksandra Faust, et al. Stop regressing: Training
324 value functions via classification for scalable deep rl. *arXiv preprint arXiv:2403.03950*, 2024.
- 325 Chelsea Finn, Xin Yu Tan, Yan Duan, Trevor Darrell, Sergey Levine, and Pieter Abbeel. Deep spatial
326 autoencoders for visuomotor learning. In *2016 IEEE International Conference on Robotics and*
327 *Automation (ICRA)*, pp. 512–519. IEEE, 2016.
- 328 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimiza-
329 tion for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- 330 Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-
331 critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- 332 Scott Fujimoto, David Meger, and Doina Precup. A deep reinforcement learning approach to
333 marginalized importance sampling with the successor representation. In *International Confer-*
334 *ence on Machine Learning*, pp. 3518–3529. PMLR, 2021.
- 335 Scott Fujimoto, Pierluca D’Oro, Amy Zhang, Yuandong Tian, and Michael Rabbat. Towards
336 general-purpose model-free reinforcement learning. *arXiv preprint arXiv:2501.16142*, 2025.

- 337 Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp:
338 Learning continuous latent space models for representation learning. In *International conference*
339 *on machine learning*, pp. 2170–2179. PMLR, 2019.
- 340 Gene H Golub and Charles F Van Loan. Lanczos methods. *Matrix Computations. Baltimore: Johns*
341 *Hopkins University Press*, pp. 470–507, 1996.
- 342 Zhaohan Guo, Shantanu Thakoor, Miruna Pîslar, Bernardo Avila Pires, Florent Alché, Corentin
343 Tallec, Alaa Saade, Daniele Calandriello, Jean-Bastien Grill, Yunhao Tang, et al. Byol-explore:
344 Exploration by bootstrapped prediction. *Advances in neural information processing systems*, 35:
345 31855–31870, 2022.
- 346 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 347 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
348 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
349 *ence on machine learning*, pp. 1861–1870. PMLR, 2018.
- 350 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
351 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 352 Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for contin-
353 uous control. *arXiv preprint arXiv:2310.16828*, 2023.
- 354 Fengxiang He, Tongliang Liu, and Dacheng Tao. Why resnet works? residuals generalize. *IEEE*
355 *transactions on neural networks and learning systems*, 31(12):5349–5362, 2020.
- 356 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
357 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
358 770–778, 2016.
- 359 Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. Pointwise convolutional neural networks. In
360 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 984–993,
361 2018.
- 362 Suning Huang, Zheyu Zhang, Tianhai Liang, Yihan Xu, Zhehao Kou, Chenhao Lu, Guowei Xu,
363 Zhengrong Xue, and Huazhe Xu. Mentor: Mixture-of-experts network with task-oriented pertur-
364 bation for visual reinforcement learning. *Proc. the International Conference on Machine Learning*
365 *(ICML)*, 2025.
- 366 Kyungsoo Kim, Jeongsoo Ha, and Yusung Kim. Self-predictive dynamics for generalization of
367 vision-based reinforcement learning. In *IJCAI*, pp. 3150–3156, 2022.
- 368 Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing
369 deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020.
- 370 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
371 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 372 Manoj Kumar, Mostafa Dehghani, and Neil Houlsby. Dual patchnorm. *arXiv preprint*
373 *arXiv:2302.01327*, 2023.
- 374 Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Rein-
375 forcement learning with augmented data. *Advances in neural information processing systems*, 33:
376 19884–19895, 2020.
- 377 Alex X Lee, Anusha Nagabandi, Pieter Abbeel, and Sergey Levine. Stochastic latent actor-critic:
378 Deep reinforcement learning with a latent variable model. *Advances in Neural Information Pro-*
379 *cessing Systems*, 33:741–752, 2020.

- 380 Hojoon Lee, Hanseul Cho, Hyunseung Kim, Daehoon Gwak, Joonkee Kim, Jaegul Choo, Se-Young
 381 Yun, and Chulhee Yun. Plastic: Improving input and label plasticity for sample efficient rein-
 382 forcement learning. *Advances in Neural Information Processing Systems*, 36, 2024a.
- 383 Hojoon Lee, Dongyoon Hwang, Donghu Kim, Hyunseung Kim, Jun Jet Tai, Kaushik Subramanian,
 384 Peter R Wurman, Jaegul Choo, Peter Stone, and Takuma Seno. Simba: Simplicity bias for scaling
 385 up parameters in deep reinforcement learning. *arXiv preprint arXiv:2410.09754*, 2024b.
- 386 Hojoon Lee, Youngdo Lee, Takuma Seno, Donghu Kim, Peter Stone, and Jaegul Choo. Hyperspher-
 387 ical normalization for scalable deep reinforcement learning. *arXiv preprint arXiv:2502.15280*,
 388 2025.
- 389 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *ArXiv e-prints*, pp.
 390 arXiv-1607, 2016.
- 391 TP Lillicrap. Continuous control with deep reinforcement learning. *arXiv preprint*
 392 *arXiv:1509.02971*, 2015.
- 393 Haotian Lin, Pengcheng Wang, Jeff Schneider, and Guanya Shi. Td-m (pc)²: Improving temporal
 394 difference mpc through policy constraint. *arXiv preprint arXiv:2502.03550*, 2025.
- 395 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.
 396 A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and*
 397 *pattern recognition*, pp. 11976–11986, 2022.
- 398 Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney.
 399 Understanding plasticity in neural networks. *Proc. the International Conference on Machine*
 400 *Learning (ICML)*, 2023.
- 401 Clare Lyle, Zeyu Zheng, Khimya Khetarpal, James Martens, Hado van Hasselt, Razvan Pascanu,
 402 and Will Dabney. Normalization and effective learning rates in reinforcement learning. *arXiv*
 403 *preprint arXiv:2407.01800*, 2024.
- 404 Guozheng Ma, Zhen Wang, Zhecheng Yuan, Xueqian Wang, Bo Yuan, and Dacheng Tao. A
 405 comprehensive survey of data augmentation in visual reinforcement learning. *arXiv preprint*
 406 *arXiv:2210.04561*, 2022.
- 407 Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin,
 408 David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance
 409 gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- 410 Trevor McInroe, Lukas Schäfer, and Stefano V Albrecht. Learning temporally-consistent represen-
 411 tations for data-efficient reinforcement learning. *arXiv preprint arXiv:2110.04935*, 2021.
- 412 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-
 413 mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
 414 control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 415 Tianwei Ni, Benjamin Eysenbach, Erfan Seyedsalehi, Michel Ma, Clement Gehring, Aditya Ma-
 416 hajan, and Pierre-Luc Bacon. Bridging state and history representations: Understanding self-
 417 predictive rl. *arXiv preprint arXiv:2401.08898*, 2024.
- 418 Johan Obando-Ceron, Ghada Sokar, Timon Willi, Clare Lyle, Jesse Farebrother, Jakob Foerster,
 419 Gintare Karolina Dziugaite, Doina Precup, and Pablo Samuel Castro. Mixtures of experts unlock
 420 parameter scaling for deep rl. *arXiv preprint arXiv:2402.08609*, 2024.
- 421 Daniel Palenicek, Florian Vogt, and Jan Peters. Scaling off-policy reinforcement learning with batch
 422 and weight normalization. *arXiv preprint arXiv:2502.07523*, 2025.

- 423 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration
424 by self-supervised prediction. In *International conference on machine learning*, pp. 2778–2787.
425 PMLR, 2017.
- 426 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
427 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
428 learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- 429 Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bach-
430 man. Data-efficient reinforcement learning with self-predictive representations. *arXiv preprint*
431 *arXiv:2007.05929*, 2020.
- 432 Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin, R De-
433 von Hjelm, Philip Bachman, and Aaron C Courville. Pretraining representations for data-efficient
434 reinforcement learning. *Advances in Neural Information Processing Systems*, 34:12686–12699,
435 2021.
- 436 Max Schwarzer, Johan Samir Obando Ceron, Aaron Courville, Marc G Bellemare, Rishabh Agar-
437 wal, and Pablo Samuel Castro. Bigger, better, faster: Human-level atari with human-level effi-
438 ciency. In *International Conference on Machine Learning*, pp. 30365–30380. PMLR, 2023.
- 439 Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak.
440 Planning to explore via self-supervised world models. In *International conference on machine*
441 *learning*, pp. 8583–8592. PMLR, 2020.
- 442 Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with
443 action-free pre-training from videos. In *International Conference on Machine Learning*, pp.
444 19561–19579. PMLR, 2022.
- 445 Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evcı. The dormant neuron phe-
446 nomenon in deep reinforcement learning. *arXiv preprint arXiv:2302.12902*, 2023a.
- 447 Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evcı. The dormant neuron phe-
448 nomenon in deep reinforcement learning. In *International Conference on Machine Learning*, pp.
449 32145–32168. PMLR, 2023b.
- 450 Ghada Sokar, Johan Obando-Ceron, Aaron Courville, Hugo Larochelle, and Pablo Samuel Cas-
451 tro. Don’t flatten, tokenize! unlocking the key to softmoe’s efficacy in deep rl. *arXiv preprint*
452 *arXiv:2410.01930*, 2024.
- 453 Adam Stooke, Kimin Lee, Pieter Abbeel, and Michael Laskin. Decoupling representation learning
454 from reinforcement learning. In *International conference on machine learning*, pp. 9870–9879.
455 PMLR, 2021.
- 456 Bhavya Sukhija, Stelian Coros, Andreas Krause, Pieter Abbeel, and Carmelo Sferrazza. Maxinforl:
457 Boosting exploration in reinforcement learning through information gain maximization. *arXiv*
458 *preprint arXiv:2412.12098*, 2024.
- 459 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud-
460 den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv*
461 *preprint arXiv:1801.00690*, 2018.
- 462 Damien Teney, Armand Mihai Nicolicioiu, Valentin Hartmann, and Ehsan Abbasnejad. Neural red-
463 shift: Random networks are not random functions. In *Proceedings of the IEEE/CVF Conference*
464 *on Computer Vision and Pattern Recognition*, pp. 4786–4796, 2024.
- 465 Raphael Trumpp, Ansgar Schäfftlein, Mirco Theile, and Marco Caccamo. Impoola: The power of
466 average pooling for image-based deep reinforcement learning. *arXiv preprint arXiv:2503.05546*,
467 2025.

- 468 Herke Van Hoof, Nutan Chen, Maximilian Karl, Patrick Van Der Smagt, and Jan Peters. Stable rein-
469 forcement learning with autoencoders for tactile and visual data. In *2016 IEEE/RSJ international*
470 *conference on intelligent robots and systems (IROS)*, pp. 3928–3934. IEEE, 2016.
- 471 Xiyao Wang, Ruijie Zheng, Yanchao Sun, Ruonan Jia, Wichayaporn Wongkamjan, Huazhe Xu, and
472 Furong Huang. Coplanner: Plan to roll out conservatively but to explore optimistically for model-
473 based rl. *arXiv preprint arXiv:2310.07220*, 2023.
- 474 Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A
475 locally linear latent dynamics model for control from raw images. *Advances in neural information*
476 *processing systems*, 28, 2015.
- 477 Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and
478 Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In
479 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16133–
480 16142, 2023.
- 481 Philipp Wu, Alejandro Escontrela, Danijar Hafner, Pieter Abbeel, and Ken Goldberg. Daydreamer:
482 World models for physical robot learning. In *Conference on robot learning*, pp. 2226–2240.
483 PMLR, 2023.
- 484 Guowei Xu, Ruijie Zheng, Yongyuan Liang, Xiyao Wang, Zhecheng Yuan, Tianying Ji, Yu Luo,
485 Xiaoyu Liu, Jiaxin Yuan, Pu Hua, et al. Drm: Mastering visual reinforcement learning through
486 dormant ratio minimization. *arXiv preprint arXiv:2310.19668*, 2023.
- 487 Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous con-
488 trol: Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021a.
- 489 Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving
490 sample efficiency in model-free reinforcement learning from images. In *Proceedings of the aaai*
491 *conference on artificial intelligence*, volume 35, pp. 10674–10681, 2021b.
- 492 Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games
493 with limited data. *Advances in neural information processing systems*, 34:25476–25488, 2021.
- 494 Tao Yu, Cuiling Lan, Wenjun Zeng, Mingxiao Feng, Zhizheng Zhang, and Zhibo Chen. Playvirtual:
495 Augmenting cycle-consistent virtual trajectories for reinforcement learning. *Advances in Neural*
496 *Information Processing Systems*, 34:5276–5289, 2021.
- 497 Tao Yu, Zhizheng Zhang, Cuiling Lan, Yan Lu, and Zhibo Chen. Mask-based latent reconstruc-
498 tion for reinforcement learning. *Advances in Neural Information Processing Systems*, 35:25117–
499 25131, 2022.
- 500 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
501 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
502 In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
- 503 Kevin Zakka, Baruch Tabanpour, Qiayuan Liao, Mustafa Haiderbhai, Samuel Holt, Jing Yuan Luo,
504 Arthur Allshire, Erik Frey, Koushil Sreenath, Lueder A Kahrs, et al. Mujoco playground. *arXiv*
505 *preprint arXiv:2502.08844*, 2025.
- 506 Ruijie Zheng, Xiyao Wang, Yanchao Sun, Shuang Ma, Jieyu Zhao, Huazhe Xu, Hal Daumé III,
507 and Furong Huang. Taco: Temporal latent action-driven contrastive loss for visual reinforcement
508 learning. *Advances in Neural Information Processing Systems*, 36:48203–48225, 2023.