
Uncovering Untapped Potential in Sample-Efficient World Model Agents

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 World model (WM) agents enable sample-efficient reinforcement learning by learn-
2 ing policies entirely from simulated experience. However, existing token-based
3 world models (TBWMs) are limited to visual inputs and discrete actions, restricting
4 their adoption and applicability. Moreover, although both intrinsic motivation and
5 prioritized WM replay have shown promise in improving WM performance and gener-
6 alization, they remain underexplored in this setting, particularly in combination.
7 We introduce Simulus, a highly modular TBWM agent that integrates (1) a modular
8 multi-modality tokenization framework, (2) intrinsic motivation, (3) prioritized
9 WM replay, and (4) regression-as-classification for reward and return prediction.
10 Simulus achieves state-of-the-art sample efficiency for planning-free WMs across
11 three diverse benchmarks. Ablation studies reveal the individual contribution of
12 each component while highlighting their synergy. Our code and model weights are
13 publicly available at <https://anonymous.4open.science/r/Simulus-FBF5>.

14 1 Introduction

15 Sample efficiency refers to the ability of a reinforcement learning (RL) algorithm to learn effective
16 policies using as few environment interactions as possible. In many real-world domains such as
17 robotics, autonomous driving, and healthcare, this is particularly critical, as interactions are costly,
18 slow, or constrained. World model agents, methods that learn control entirely from simulated
19 experience generated by a learned dynamics model, have emerged as a promising approach to
20 improving sample efficiency [17, 20, 37, 57].

21 Here, we focus on token-based world models (TBWMs) [37, 38, 11], sample-efficient RL methods
22 that learn the dynamics entirely within a learned discrete token space, where each observation
23 comprises a sequence of tokens. Evidently, most large scale world models [1, 12, 14] operate on
24 multi-token observations, suggesting that such representations are advantageous at scale. TBWMs
25 offer a clear modular design, separating the optimization of its representation, dynamics, and control
26 models. As modular systems, TBWMs are easier to scale, develop, study, and deploy, as individual
27 modules can be treated independently without interfering and are easier to master through divide and
28 conquer. In addition, such separation leads to simpler optimization objectives and avoids interference
29 between them (Appendix B.1).

30 However, existing TBWMs are restricted to image observations and discrete actions, such as Atari
31 games, limiting both their adoption and broader applicability, as their effectiveness in diverse
32 environments and modalities remains unclear. While multi-modality tokenization approaches exist for
33 large-scale offline settings [43, 47, 33, 34], these methods rely on large vocabularies (e.g., 33K tokens)
34 which are inefficient for online, data-limited regimes. Whether substantially smaller vocabularies can
35 preserve competitive precision and performance is still an open question, leaving these approaches
36 underexplored for sample-efficient RL.

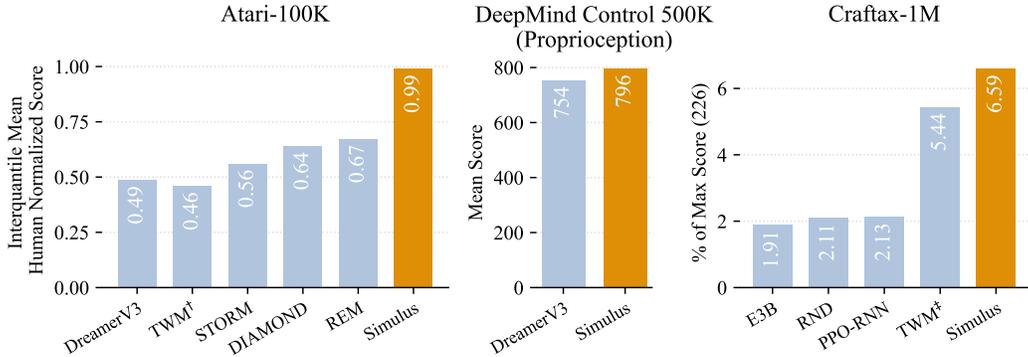


Figure 1: Results overview. Simulus exhibits state-of-the-art sample-efficiency performance for planning-free methods across all three benchmarks. † [44], ‡ [13].

37 Furthermore, despite compelling results [46, 49, 30], intrinsic motivation and prioritized *world model*
 38 replay remain underexplored in sample-efficient world model agents [19, 20, 60, 11, 3], particularly
 39 in combination. We conjecture that intrinsic motivation is underused as it may steer the agent toward
 40 task-irrelevant regions, potentially wasting limited interaction budget. Prioritized replay [30], while
 41 promising, lacks robust empirical support and proved challenging to tune in our experiments.

42 To address these limitations, we propose Simulus, a modular world model that extends a recent
 43 TBWM method [11] by integrating several powerful advances from the literature: (1) a modular
 44 multi-modality tokenization framework for handling arbitrary combinations of observation and action
 45 modalities, (2) intrinsic motivation for epistemic uncertainty reduction [50, 49], (3) prioritized world
 46 model replay [30], and (4) regression-as-classification (RaC) for reward and return prediction [16, 20].

47 To evaluate the impact of the proposed components, we conducted extensive empirical evaluations
 48 across three diverse benchmarks, ranging from the visual Atari 100K [28], to the continuous proprio-
 49 ception tasks of the DeepMind Control Suite [53], to Craftax [35], which combines symbolic 2D grid
 50 maps with continuous state features. There, Simulus achieves state-of-the-art sample-efficiency for
 51 planning-free world models. Ablation studies further show the contribution and effectiveness of each
 52 component.

53 Summary of contributions:

- 54 • Through extensive empirical evaluations across three diverse benchmarks, we show that
 55 intrinsic motivation, prioritized replay, and RaC significantly improve performance in the
 56 sample-efficiency setting for world model agents, particularly when combined.
- 57 • Our results demonstrate the effectiveness of the multi-modality tokenization approach, as
 58 Simulus achieves state-of-the-art planning-free performance across benchmarks, establishing
 59 TBWMs as widely-applicable methods.
- 60 • We propose Simulus, a versatile TBWM agent that follows a highly modular design, offering
 61 a solid foundation for future developments. To support future research and adoption of
 62 TBWMs, we open-source our code and weights.

63 2 Method

64 **Notations** We consider the Partially Observable Markov Decision Process (POMDP) setting. How-
 65 ever, since in practice the agent has no knowledge about the hidden state space, consider the following
 66 state-agnostic formulation. Let Ω, \mathcal{A} be the sets of observations and actions, respectively. At every
 67 step t , the agent observes $\mathbf{o}_t \in \Omega$ and picks an action $\mathbf{a}_t \in \mathcal{A}$. From the agent’s perspective, the envi-
 68 ronment evolves according to $\mathbf{o}_{t+1}, r_t, d_t \sim p(\mathbf{o}_{t+1}, r_t, d_t | \mathbf{o}_{\leq t}, \mathbf{a}_{\leq t})$, where r_t, d_t are the observed
 69 reward and termination signals, respectively. The process repeats until a positive termination signal
 70 $d_t \in \{0, 1\}$ is obtained. The agent’s objective is to maximize its expected return $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_{t+1}]$
 71 where $\gamma \in [0, 1]$ is a discount factor.

72 For multi-modal observations, let $\mathbf{o}_t = \{\mathbf{o}_t^{(i)}\}_{i=1}^{|\kappa|}$ where κ is the set of environment modalities and
 73 $\mathbf{o}_t^{(i)}$ denotes the features of modality κ_i .

74 **Overview** Simulus builds on REM [11]. The agent comprises a representation module \mathcal{V} , a world
 75 model \mathcal{M} , a controller \mathcal{C} , and a replay buffer. To facilitate a modular design, following REM, each
 76 module is optimized separately. The training process of the agent involves a repeated cycle of four
 77 steps: data collection, representation learning (\mathcal{V}), world model learning (\mathcal{M}), and control learning in
 78 imagination (\mathcal{C}).

79 2.1 The Representation Module \mathcal{V}

80 \mathcal{V} is responsible for encoding and decoding raw observations and actions. It is a modular tokenization
 81 system with encoder-decoder pairs for different input modalities. Encoders produce fixed-length
 82 token sequences, creating a common interface that enables combining tokens from various sources
 83 into a unified representation. After embedding, these token sequences are concatenated into a single
 84 representation, as described in Section 2.2. Note that encoder-decoder pairs need not be learning-
 85 based methods; however, when learned, they are optimized independently. This design enables \mathcal{V} to
 86 deal with any combination of input modalities, provided the respective encoder-decoder pairs.

87 **Tokenization** \mathcal{V} transforms raw observations \mathbf{o} to sets of fixed-length integer token sequences
 88 $\mathbf{z} = \{\mathbf{z}^{(i)}\}_{i=1}^{|\kappa|}$ by applying the encoder of each modality $\mathbf{z}^{(i)} = \text{enc}_i(\mathbf{o}^{(i)})$. Actions \mathbf{a} are tokenized
 89 using the encoder-decoder pair of the related modality to produce \mathbf{z}^a . The respective decoders
 90 reconstruct observations from their tokens: $\hat{\mathbf{o}}^{(i)} = \text{dec}_i(\mathbf{z}^{(i)})$.

91 Simulus natively supports four modalities: images, continuous vectors, categorical variables, and
 92 image-like multi-channel grids of categorical variables, referred to as "2D categoricals". More
 93 formally, 2D categoricals are elements of $([k_1] \times [k_2] \times \dots \times [k_C])^{m \times n}$ where k_1, \dots, k_C are per
 94 channel vocabulary sizes, C is the number of channels, m, n are spatial dimensions, and $[k] =$
 95 $\{1, \dots, k\}$.

96 Following REM, we use a VQ-VAE [15, 55] for image observations. For the tokenization of
 97 continuous vectors, each feature is quantized to produce a token, as in [43]. However, to improve
 98 learning efficiency, we reduce the vocabulary size from 1024 to 125, and modify the quantization
 99 levels for optimal coverage (Appendix A.2.2). Unbounded vectors are first transformed using the
 100 symlog function [20], defined as $\text{symlog}(x) = \text{sign}(x) \ln(1 + |x|)$, which compresses the magnitude
 101 of large absolute values. Lastly, while no special tokenization is required for categorical inputs, 2D
 102 categoricals are flattened along the spatial dimensions to form a sequence of categorical vectors. The
 103 embedding of each token vector is obtained by averaging the embeddings of its entries.

104 2.2 The World Model \mathcal{M}

105 The purpose of \mathcal{M} is to learn a model of the environment’s dynamics. Concretely, given trajectory
 106 segments $\tau_t = \mathbf{z}_1, \mathbf{z}_1^a, \dots, \mathbf{z}_t, \mathbf{z}_t^a$ in token representation, \mathcal{M} models the distributions of the next
 107 observation and termination signal, and the expected reward:

$$\text{Transition: } p_\theta(\hat{\mathbf{z}}_{t+1} | \tau_t), \quad (1)$$

$$\text{Reward: } \hat{r}_t = \hat{r}_\theta(\tau_t), \quad (2)$$

$$\text{Termination: } p_\theta(\hat{d}_t | \tau_t), \quad (3)$$

108 where θ is the parameters vector of \mathcal{M} and $\hat{r}_\theta(\tau_t)$ is an estimator of $\mathbb{E}_{r_t \sim p(r_t | \tau_t)}[r_t]$.

109 **Architecture** \mathcal{M} comprises a sequence model f_θ and multiple heads for the prediction of tokens
 110 of different observation modalities, rewards, termination signals, and for the estimation of model
 111 uncertainty. Concretely, f_θ is a retentive network (RetNet) [51] augmented with a parallel observation
 112 prediction (POP) [11] mechanism. All heads are implemented as multilayer perceptrons (MLP) with
 113 a single hidden layer. We defer the details about these architectures to Appendix A.3.

114 **Embedding** \mathcal{M} translates token trajectories τ into sequences of d -dimensional embeddings \mathbf{X}
 115 using a set of embedding (look-up) tables. By design, each modality is associated with a separate

116 table. In cases where an embedding table is not provided by the appropriate encoder-decoder pair, \mathcal{M}
 117 and \mathcal{C} learn dedicated tables separately and independently. As embeddings sequences are composed
 118 hierarchically, we use the following hierarchical notation:

$$\begin{aligned} \text{Observation-action block: } \mathbf{X}_t &= (\mathbf{X}_t^o, \mathbf{X}_t^a), \\ \text{Observation block: } \mathbf{X}_t^o &= (\mathbf{X}_t^{(1)}, \dots, \mathbf{X}_t^{(|\kappa|)}), \end{aligned}$$

119 where K_i denotes the number of embedding vectors in $\mathbf{X}_t^{(i)}$. Similarly, $K = \sum_{i=1}^{|\kappa|} K_i$. To combine
 120 latents of each \mathbf{z}_t , \mathcal{V} concatenates their token sequences along the temporal axis based on a predefined
 121 modality order. We defer the full details on the embedding process to Appendix A.3.

122 **Sequence Modeling** Given a sequence of
 123 observation-action blocks $\mathbf{X} = \mathbf{X}_1, \dots, \mathbf{X}_t$,
 124 the matching outputs $\mathbf{Y}_1, \dots, \mathbf{Y}_t$ are com-
 125 puted auto-regressively as follows:

$$(\mathbf{S}_t, \mathbf{Y}_t) = f_\theta(\mathbf{S}_{t-1}, \mathbf{X}_t),$$

126 where \mathbf{S}_t is a recurrent state that summarizes
 127 $\mathbf{X}_{\leq t}$ and $\mathbf{S}_0 = 0$. However, the output \mathbf{Y}_{t+1}^u ,
 128 from which $\hat{\mathbf{z}}_{t+1}$ is predicted, is computed
 129 using the POP mechanism via another call as

$$(\cdot, \mathbf{Y}_{t+1}^u) = f_\theta(\mathbf{S}_t, \mathbf{X}^u),$$

130 where $\mathbf{X}^u \in \mathbb{R}^{K \times d}$ is a learned embedding se-
 131 quence. Intuitively, \mathbf{X}^u acts as a learned prior,
 132 enabling the parallel generation of multiple
 133 tokens into the future.

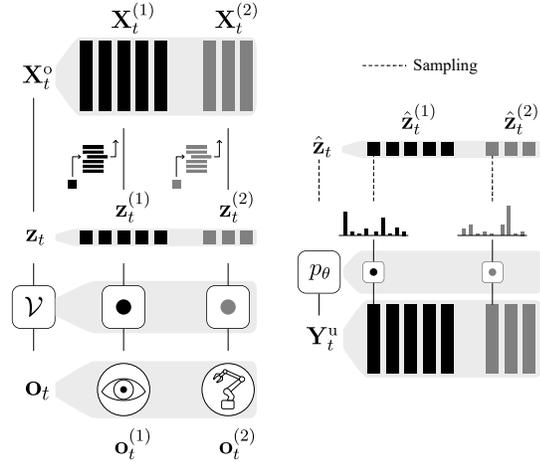
134 To model $p_\theta(\hat{\mathbf{z}}_{t+1} | \mathbf{Y}_{t+1}^u)$, the distributions
 135 $p_\theta(\hat{z} | \mathbf{y})$ of each token \hat{z} of each modal-
 136 ity κ_i are modeled using modality-specific pre-
 137 diction heads implemented as MLPs with a
 138 single hidden layer and an output size equal
 139 to the vocabulary size of enc_i (Figure 2b). For
 140 2D categoricals, C heads are used to predict
 141 the C tokens from each \mathbf{y} .

142 Similarly, rewards and termination signals are predicted by additional prediction heads as $\hat{r}_t =$
 143 $\hat{r}_\theta(\mathbf{y})$, $\hat{d}_t \sim p_\theta(\hat{d}_t | \mathbf{y})$, slightly abusing notations, where \mathbf{y} is the last vector of \mathbf{Y}_{t+1}^u . An illustration
 144 is provided in Figure 3.

145 **Reducing Epistemic Uncertainty via Intrinsic Motivation** The world model \mathcal{M} serves as the
 146 cornerstone of the entire system. Any controller operating within a world model framework can only
 147 perform as well as the underlying world model allows, making its quality a fundamental limiting
 148 factor. In deep learning methods, the model’s performance depends heavily on the quality of its
 149 training data. Accurate dynamics modeling requires comprehensive data collection that captures the
 150 full spectrum of possible environmental behaviors. This presents a particular challenge in online RL,
 151 where the controller must systematically and efficiently explore its environment. Success depends on
 152 intelligently guiding the controller toward unexplored or undersampled regions of the dynamics space.
 153 An effective approach to this challenge involves estimating the world model’s epistemic uncertainty
 154 and directing the controller to gather data from regions where this uncertainty is highest [46, 50, 49].

155 Our approach estimates epistemic uncertainty using an ensemble of $N_{\text{ens}} = 4$ next observation
 156 prediction heads $\{p_{\phi_i}(\hat{\mathbf{z}} | \text{sg}(\mathbf{Y}^u))\}_{i=1}^{N_{\text{ens}}}$ with parameters $\{\phi_i\}_{i=1}^{N_{\text{ens}}}$ [49, 31] where $\text{sg}(\cdot)$ is the stop
 157 gradient operator. To quantify disagreement between the ensemble’s distributions, we employ the
 158 Jensen-Shannon divergence (JSD) [50]. For probability distributions P_1, \dots, P_n , the JSD is defined
 159 as:

$$\text{JSD}(P_1, \dots, P_n) = \mathcal{H}\left(\frac{1}{n} \sum_{i=1}^n P_i\right) - \frac{1}{n} \sum_{i=1}^n \mathcal{H}(P_i),$$



(a) Observation tokeniza- (b) Prediction of observa-
 tion and embedding. tion tokens.

Figure 2: An illustration of the independent pro-
 cessing of modalities for an observation with two
 modalities.

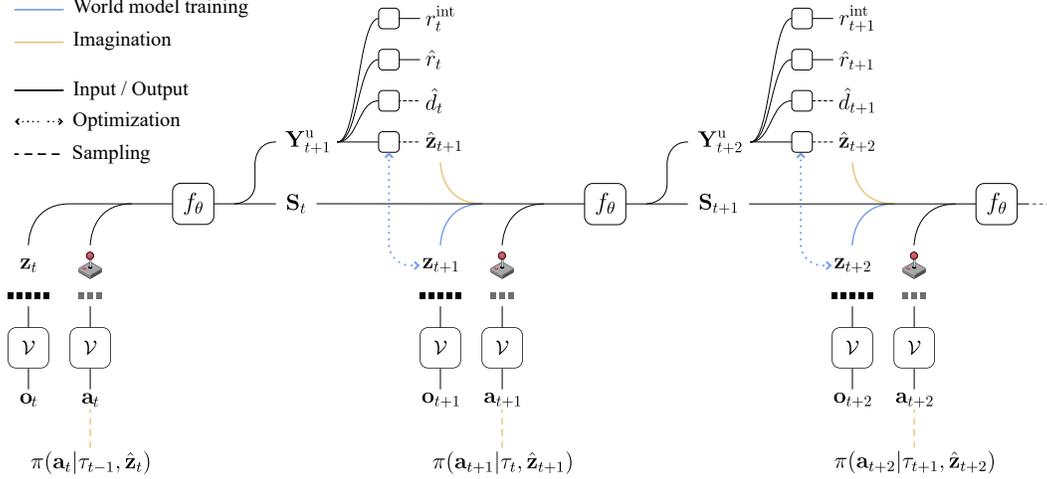


Figure 3: World model training and imagination. To maintain visual clarity, we omitted token embedding details, as well as optimization details of rewards and termination signals.

160 where $\mathcal{H}(\cdot)$ denotes the Shannon entropy. Since observations comprise multiple tokens, we average
 161 the per-token JSD values to obtain a single uncertainty measure δ_t . Training data is divided equally
 162 among ensemble members, with each predictor processing a distinct subset of each batch. Despite
 163 the ensemble approach, our implementation maintains computational efficiency, with negligible
 164 additional overhead in practice.

165 To guide \mathcal{C} towards regions of high epistemic uncertainty, \mathcal{M} provides \mathcal{C} with additional intrinsic
 166 rewards $r_t^{\text{int}} = \delta_t$ during imagination. Here, the reward provided by \mathcal{M} at each step t is given by

$$\bar{r}_t = w^{\text{int}} r_t^{\text{int}} + w^{\text{ext}} \hat{r}_t,$$

167 where $w^{\text{int}}, w^{\text{ext}} \in \mathbb{R}$ are hyperparameters that control the scale of each reward type. Optimizing
 168 the controller in imagination allows it to reach areas of high model uncertainty without additional
 169 real-environment interaction.

170 **Prioritized Replay** Recent work has demonstrated that prioritizing replay buffer sampling during
 171 world model training could lead to significant performance gains in intrinsically motivated agents
 172 [30]. While their approach showed promise, it required extensive hyperparameter tuning in practice.
 173 We propose a simpler, more robust prioritization scheme for world model training.

174 Here, the replay buffer maintains a world model loss value for each stored example, with newly
 175 added examples assigned a high initial loss value ν_0 . During \mathcal{M} 's training, we sample each batch
 176 using a mixture of uniform and prioritized sampling, controlled by a single parameter $\alpha \in [0, 1]$
 177 that determines the fraction of prioritized samples. For the prioritized portion, we sample examples
 178 proportional to their softmax-transformed losses $p_i = \text{softmax}(\mathcal{L})_i$. The loss values are updated
 179 after each world model optimization step using the examples' current batch losses.

180 **Training** We use the cross-entropy loss for the optimization of all components of \mathcal{M} . Specifically,
 181 for each t , the loss of $p_\theta(\hat{\mathbf{z}}_t | \mathbf{Y}_t^u)$ is obtained by averaging the cross-entropy losses of its individual
 182 tokens. The same loss is used for each ensemble member $p_{\phi_i}(\hat{\mathbf{z}}_t | \text{sg}(\mathbf{Y}_t^u))$. The optimization and
 183 design of the reward predictor is similar to that of the critic, as described in Section 2.3. A formal
 184 description of the optimization objective can be found in Appendix A.3.1.

185 2.3 The Controller \mathcal{C}

186 \mathcal{C} is an extended version of the actor-critic of REM [11] that supports additional observation and
 187 action spaces and implements regression-as-classification for return predictions.

188 **Architecture** At the core of \mathcal{C} 's architecture, parameterized by ψ , is an LSTM [26] sequence model.
 189 At each step t , upon observing \mathbf{z}_t , a set of modality-specific encoders map each modality tokens $\mathbf{z}_t^{(i)}$ to

190 a latent vector $\mathbf{x}^{(i)}$, where we abuse our notation \mathbf{x} as the context of the discussion is limited to \mathcal{C} . The
 191 latents are then fused by a fully-connected network to obtain a single vector $\mathbf{x} = g_\psi(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(|\kappa|)})$.
 192 $\mathbf{x}_t \in \mathbb{R}^{d_c}$ is processed by \mathcal{C} 's sequence model to produce $\mathbf{h}_t, \mathbf{c}_t = \text{LSTM}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}; \psi)$ where
 193 $\mathbf{h}_t, \mathbf{c}_t$ are the LSTM's hidden and cell states, respectively. Lastly, two linear output layers produce the
 194 logits from which the actor and critic outputs $\pi(\mathbf{a}_t|\mathbf{h}_t), \hat{V}^\pi(\mathbf{h}_t)$ are derived. For continuous action
 195 spaces, the actor uses a categorical distribution over a uniformly spaced discrete subset of $[-1, 1]$.
 196 We defer the full details about the encoding process to Appendix A.4.

197 **Regression as Classification for Reward and Return Prediction** Robustly handling unbounded
 198 reward signals has long been challenging as they can vary dramatically in both magnitude and
 199 frequency. Hafner et al. [20] addressed this challenge by using a classification network that predicts
 200 the weights of exponentially spaced bins and employed a two-hot loss for the network's optimization.
 201 Farebrother et al. [16] studied the use of cross-entropy loss in place of the traditional mean squared
 202 error loss for value-based deep RL methods. In their work, the HL-Gauss method was shown
 203 to significantly outperform the two-hot loss method. Building on these developments, we adopt
 204 the classification network with exponential bins from [20], and apply the HL-Gauss method for
 205 its optimization. Concretely, the critic's value estimates are predicted using a linear output layer
 206 parameterized by $\mathbf{W} \in \mathbb{R}^{m \times d_c}$ with $m = 128$ outputs corresponding to m uniform bins defined by
 207 $m + 1$ endpoints $\mathbf{b} = (b_0, \dots, b_m)$. The predicted value is given by

$$\hat{y} = \text{symexp}\left(\text{softmax}(\mathbf{W}\mathbf{h})^\top \hat{\mathbf{b}}\right)$$

208 where $\text{symexp}(x) = \text{sign}(x)(\exp(|x|) - 1)$ is the inverse of the symlog function and $\hat{\mathbf{b}} =$
 209 $\left(\frac{b_1+b_0}{2}, \dots, \frac{b_m+b_{m-1}}{2}\right)$ are the bin centers. Given the true target $y \in \mathbb{R}$, the HL-Gauss loss is
 210 given by

$$\mathcal{L}_{\text{HL-Gauss}}(\mathbf{W}, \mathbf{h}, y) = \tilde{\mathbf{y}}^\top \log \text{softmax}(\mathbf{W}\mathbf{h})$$

211 where $\tilde{y}_i = \Phi\left(\frac{b_i - \text{symlog}(y)}{\sigma}\right) - \Phi\left(\frac{b_{i-1} - \text{symlog}(y)}{\sigma}\right)$, Φ is the cumulative density function of the
 212 standard normal distribution and σ is a standard deviation hyperparameter that controls the amount of
 213 label smoothing.

214 **Training in Imagination** \mathcal{C} is trained entirely from simulated experience generated through interac-
 215 tion with \mathcal{M} . Specifically, \mathcal{M} and \mathcal{C} are initialized with a short trajectory segment sampled uniformly
 216 from the replay buffer and interact for H steps. An illustration of this process is given in Figure 3
 217 (orange path). λ -returns are computed for each generated trajectory segment and are used as targets
 218 for critic learning. For policy learning, a REINFORCE [52] objective is used, with a \hat{V}^π baseline for
 219 variance reduction. See Appendix A.4.2 for further details.

220 3 Experiments

221 To evaluate sample efficiency, we used benchmarks that measure performance within a fixed, limited
 222 environment interaction budget. These benchmarks were also selected to address key research
 223 questions: (1) whether the proposed multi-modality approach is effective—both in continuous control
 224 settings and in handling multi-modal observations; and (2) whether the integrated components are
 225 effective across diverse environments (Section 3.3).

226 3.1 Experimental Setup

227 **Benchmarks:** We evaluate Simulus on three sample-efficiency benchmarks of different observation
 228 and action modalities: Atari 100K [28], DeepMind Control Suite (DMC) Proprioception 500K [53],
 229 and Craftax-1M [35].

230 Atari 100K has become the gold standard in the literature for evaluating sample-efficient deep RL
 231 agents. The benchmark comprises a subset of 26 games. Within each game, agents must learn from
 232 visual image signal under a tightly restricted budget of 100K interactions, corresponding to roughly
 233 two hours of human gameplay.

234 The DeepMind Control Suite (DMC) is a set of continuous control tasks involving multiple agent
 235 embodiments ranging from simple single-joint models to complex humanoids. Here, we follow the

subset of proprioception tasks used for the evaluation of DreamerV3 [20], where observations and actions are continuous vectors. At each task, the agent’s interaction budget is limited to 500K steps.

Craftax is a 2D open-world survival game benchmark inspired by Minecraft, designed to evaluate RL agents’ capabilities in planning, memory, and exploration. The partially-observable environment features procedurally generated worlds where agents must gather and craft resources while surviving against hostile creatures. Observations consist of a 9×11 tile egocentric map, where each tile consists of 4 symbols, and 48 state features corresponding to state information such as inventory and health. Here, we consider the sample-efficiency oriented Craftax-1M variant which only allows an interaction budget of one million steps.

Baselines On Atari-100K, we compare Simulus against DreamerV3 [20] and several methods restricted to image observations: TWM [44], STORM [60], DIAMOND [3], and REM [11]. On DMC, we compare exclusively with DreamerV3, currently the only planning-free world model method with published results on the 500K proprioception benchmark. On Craftax-1M, we compare against TWM [13], a concurrent work that proposes a Transformer based world model with a focus on the Craftax benchmark, and the baselines reported in the Craftax paper: Random Network Distillation (RND) [10], PPO [48] with a recurrent neural network (PPO-RNN), and Exploration via Elliptical Episodic Bonuses (E3B) [23]. As Craftax is a recent benchmark, there are no other published results in existing world models literature. Following the standard practice in the literature, we exclude planning-based methods [21, 57], as planning is an orthogonal component that operates on any given model, typically incurring significant computational overhead.

Metrics and Evaluation For Atari, we report human-normalized scores (HNS), calculated as $\frac{\text{agent_score} - \text{random_score}}{\text{human_score} - \text{random_score}}$ [39]. Following the protocol of Agarwal et al. [2] and using their toolkit, we report the mean, median, interquartile mean (IQM), and optimality gap metrics with 95% stratified bootstrap confidence intervals. For DMC and Craftax, we report the raw agent returns. We use 5 random seeds per environment. In each experiment, final performance is evaluated using 100 test episodes at the end of training and the mean score is reported.

3.2 Results

Simulus achieves state-of-the-art performance across all three benchmarks (Figure 1). On Atari 100k, it outperforms all baselines across key metrics (Figure 4). Notably, Simulus is the first planning-free world model to reach human-level IQM and median scores, achieving superhuman performance

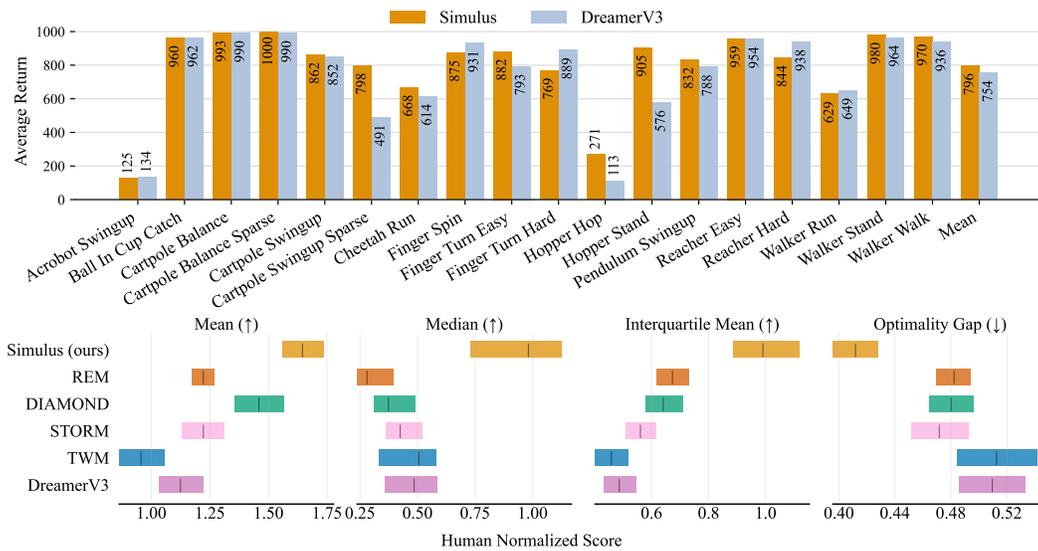


Figure 4: Results on the DeepMind Control Suite 500K Proprioception (top) and Atari 100K (bottom) benchmarks.

266 on **13** out of 26 games (Table 9, Appendix B). Building on REM, these gains are attributed to the
 267 integration of the proposed components, demonstrating their combined effectiveness.

268 **Effectiveness in continuous environments** Figure 4 provides
 269 compelling evidence that token-based architectures can perform
 270 well in continuous domains—even with compact vocabularies:
 271 Simulus consistently matches DreamerV3 across most tasks and
 272 slightly outperforms it on average.

273 **Effectiveness in environments with multi-modal observations**
 274 We evaluate multi-modality performance in Crafttax, as it com-
 275 bines an image-like 2D grid map with a vector of features, involv-
 276 ing multiple tokenizers (\mathcal{V}). Simulus maintains sample-efficiency
 277 in this multi-modal environment, outperforming both concurrent
 278 world model methods (Figure 1) and all model-free baselines
 279 (Figure 5), including exploration-focused algorithms. With 444
 280 tokens per observation arranged into sequences of 147 embed-
 281 dings, even short trajectories in Crafttax contain thousands of
 282 tokens, demonstrating Simulus’s efficient handling of long se-
 283 quences. These findings indicate that the world model (\mathcal{M}) and
 284 controller (\mathcal{C}) maintain strong performance under multi-modal
 285 inputs when processed by the proposed modular tokenizer.

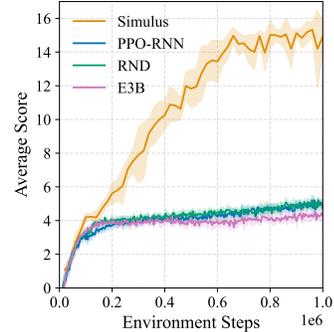


Figure 5: Crafttax-1M training curves with mean and 95% confidence intervals.

286 3.3 Ablation Studies

287 We ablate the intrinsic rewards, prioritized replay, and regression-as-classification to demon-
 288 strate their individual contributions to Simulus’s performance. In each experiment, Simulus is modified
 289 by disabling a single component. Due to limited computational resources, we consider a subset of
 290 8 tasks for each of the Atari and DMC benchmarks, and exclude Crafttax from this analysis. For
 291 Atari 100K, we used games in which significant improvements were observed. For DMC, we chose a
 292 subset that includes different embodiments. We defer the specific environment names to Appendix C.

293 The results are presented in Figure 6. Although all components contributed to Simulus’s final
 294 performance, intrinsic rewards were especially crucial for achieving competitive performance in both
 295 benchmarks. Interestingly, the Atari 100K results indicate that combining all three components yields
 296 a significantly stronger algorithm. These findings also suggest that both prioritized world model
 297 replay and regression-as-classification enhance the effectiveness of intrinsic rewards.

298 More broadly, the results in Figure 6 demonstrate that encouraging the controller to explore regions of
 299 high epistemic uncertainty through intrinsic rewards significantly improves its performance in world
 300 model agents, even in reward-rich environments. This observation is non-trivial in a sample-efficient
 301 setting, where the limited interaction budget makes model-driven exploration particularly costly, as it
 302 consumes resources that could otherwise be used for task-related exploration during data collection.
 303 The latter type of exploration aims to collect new information about the true reward signal, which
 304 defines the task and its success metric. On the other hand, model-driven exploration may guide the
 305 controller towards environment regions that are irrelevant to the task at hand.

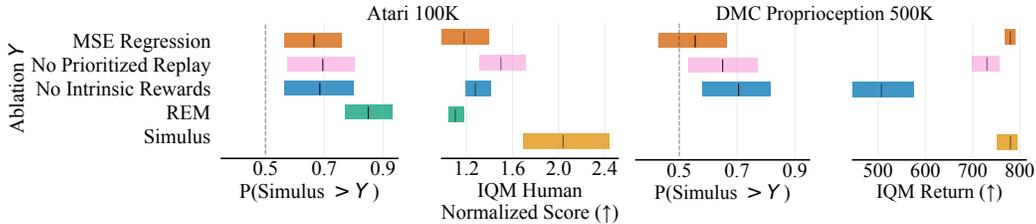


Figure 6: Ablations results on the Atari-100K and DeepMind Control Proprioception 500K bench-
 marks. A subset of 8 games was used for each ablation.

306 4 Related Work

307 **Offline Multi-Modal Methods** Large-scale token sequence models for multi-modal agent trajectories have been proposed in [33, 34, 43, 47]. Gato [43] and TDM [47] tokenize inputs via predefined transformations, while Unified IO [33, 34] leverages pretrained models. These methods do not
308 learn control through RL but rely on expert data. They also use massive models—with billions of
309 parameters, large vocabularies, and significantly more data and compute than sample-efficient world
310 models. Consequently, it remains unclear whether their design choices would be effective in online,
311 sample-efficient settings with non-stationary and limited data.
312
313

314 **World Model Agents** Model-based agents that learn policies solely from simulated data generated
315 by a learned world model were introduced by Ha and Schmidhuber [17], followed by the influential
316 Dreamer family [18–20]. Dreamer jointly optimizes its representation and recurrent world models via
317 a KL divergence between learned prior and posterior distributions, leading to interfering objectives
318 (Appendix B.1) and a complex, monolithic architecture that complicates development and scaling [59].
319 With the rise of Transformer architectures in language modeling [56, 9], Transformer-based Dreamer
320 variants emerged [60, 44], alongside token-based world models (TBWMs) that treat trajectories as
321 discrete token sequences [37, 11]. However, these methods are limited to visual environments with
322 discrete actions (e.g., Atari 100K), leaving their performance in other modalities uncertain. Recently,
323 DIAMOND [3], a diffusion world model inspired by advances in generative modeling [45], was
324 introduced. While it generates visually compelling outputs, it remains limited to visual domains.

325 **Intrinsically Motivated World Model Agents** Although intrinsic motivation (IM) has been extensively
326 studied [41, 10, 22, 5, 4], its use in world models typically involves an exploration pretraining
327 phase followed by limited task-specific fine-tuning [49, 36, 30]. While combining IM with prioritized
328 replay has shown promise [30], it remains unexplored in standard sample-efficiency settings with
329 external rewards.

330 **Large-Scale Video World Models** Building on recent advances in video generative modeling
331 [25, 7, 8], recent works have introduced large-scale video world models [14, 1, 12, 54], trained offline
332 on extensive pre-collected data to predict future frames. However, these models do not address control
333 learning, particularly RL. While recent efforts aim to bridge this gap [59], they remain confined to
334 visual environments and lack comprehensive empirical evaluation.

335 5 Limitations and Future Work

336 Here, we briefly highlight several limitations of this work. First, although the feature quantization
337 approach for tokenizing continuous vectors showed promise, it leads to excessive sequence lengths.
338 We believe that more efficient solutions can be found for dealing with continuous inputs. Second, due
339 to the scarcity of rich multi-modal RL benchmarks, we could not extensively explore diverse modality
340 combinations in our experiments. Lastly, token-based world model agents remain significantly slower
341 to train than other baselines in sample-efficient RL. Nonetheless, their modular design enables faster
342 policy inference as the controller is independent of the world model.

343 6 Conclusions

344 In this paper, we demonstrated the effectiveness of several underutilized techniques for improving
345 world model agents. A modular multi-modality tokenization framework broadens their applicability
346 across diverse domains, while intrinsic motivation, prioritized world model replay, and regression-
347 as-classification enhance sample efficiency, particularly when combined. These techniques were
348 incorporated into a token-based world model, yielding Simulus. Extensive experiments show that
349 Simulus achieves state-of-the-art performance across diverse benchmarks, including visual, contin-
350 uous, and symbolic domains. It outperforms all baselines on key metrics in both Atari-100K and
351 the challenging Craftax benchmarks. Ablations further highlight the individual contribution of each
352 component. We hope that the highly modular design of Simulus, along with the released code and
353 model weights, provides a strong foundation for future work.

354 **References**

- 355 [1] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit
356 Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model
357 platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025.
- 358 [2] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Belle-
359 mare. Deep reinforcement learning at the edge of the statistical precipice. In M. Ranzato,
360 A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in
361 Neural Information Processing Systems*, volume 34, pages 29304–29320. Curran Associates,
362 Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/file/
363 f514cec81cb148559cf475e7426eed5e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/f514cec81cb148559cf475e7426eed5e-Paper.pdf).
- 364 [3] Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and
365 François Fleuret. Diffusion for world modeling: Visual details matter in atari. *arXiv preprint
366 arXiv:2405.12399*, 2024.
- 367 [4] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvit-
368 skyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the Atari human
369 benchmark. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International
370 Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*,
371 pages 507–517. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.press/v119/
372 badia20a.html](https://proceedings.mlr.press/v119/badia20a.html).
- 373 [5] Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, Bilal Piot, Steven
374 Kapturowski, Olivier Tieleman, Martin Arjovsky, Alexander Pritzel, Andrew Bolt, and Charles
375 Blundell. Never give up: Learning directed exploration strategies. In *International Con-
376 ference on Learning Representations*, 2020. URL [https://openreview.net/forum?id=
377 Sye57xStvB](https://openreview.net/forum?id=Sye57xStvB).
- 378 [6] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients
379 through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- 380 [7] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Do-
381 minik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion:
382 Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 383 [8] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr,
384 Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh.
385 Video generation models as world simulators. 2024. URL [https://openai.com/research/
386 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 387 [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
388 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
389 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
390 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
391 teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish,
392 Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners.
393 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in
394 Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates,
395 Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/
396 1457c0d6bfbcb4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfbcb4967418bfb8ac142f64a-Paper.pdf).
- 397 [10] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random
398 network distillation. In *International Conference on Learning Representations*, 2019. URL
399 <https://openreview.net/forum?id=H1lJJnR5Ym>.
- 400 [11] Lior Cohen, Kaixin Wang, Bingyi Kang, and Shie Mannor. Improving token-based world
401 models with parallel observation prediction. In *Forty-first International Conference on Machine
402 Learning*, 2024. URL <https://openreview.net/forum?id=Lfp5Dk1xb6>.
- 403 [12] Decart, Julian Quevedo, Quinn McIntyre, Spruce Campbell, Xinlei Chen, and Robert Wachen.
404 Oasis: A universe in a transformer, 2024. URL <https://oasis-model.github.io/>.

- 405 [13] Antoine Dedieu, Joseph Ortiz, Xinghua Lou, Carter Wendelken, Wolfgang Lehrach, J Swaroop
406 Guntupalli, Miguel Lazaro-Gredilla, and Kevin Patrick Murphy. Improving transformer world
407 models for data-efficient rl, 2025. URL <https://arxiv.org/abs/2502.01591>.
- 408 [14] Google DeepMind. Genie 2: A large-scale foundation world model, 2024. URL
409 [https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-
410 world-model/](https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/).
- 411 [15] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution
412 image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
413 recognition*, pages 12873–12883, 2021.
- 414 [16] Jesse Farebrother, Jordi Orbay, Quan Vuong, Adrien Ali Taiga, Yevgen Chebotar, Ted Xiao,
415 Alex Irpan, Sergey Levine, Pablo Samuel Castro, Aleksandra Faust, Aviral Kumar, and Rishabh
416 Agarwal. Stop regressing: Training value functions via classification for scalable deep RL. In
417 *Forty-first International Conference on Machine Learning*, 2024. URL [https://openreview.
418 net/forum?id=dVpFKfqF3R](https://openreview.net/forum?id=dVpFKfqF3R).
- 419 [17] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In
420 *Advances in Neural Information Processing Systems 31*, pages 2451–2463. Curran Associates,
421 Inc., 2018. URL [https://papers.nips.cc/paper/7512-recurrent-world-models-
422 facilitate-policy-evolution](https://papers.nips.cc/paper/7512-recurrent-world-models-facilitate-policy-evolution). <https://worldmodels.github.io>.
- 423 [18] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control:
424 Learning behaviors by latent imagination. In *International Conference on Learning Representations*,
425 2020. URL <https://openreview.net/forum?id=S110TC4tDS>.
- 426 [19] Danijar Hafner, Timothy P. Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari
427 with discrete world models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=0oabwyZb0u>.
- 430 [20] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
431 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 432 [21] Nicklas Hansen, Hao Su, and Xiaolong Wang. TD-MPC2: Scalable, robust world models for
433 continuous control. In *The Twelfth International Conference on Learning Representations*, 2024.
434 URL <https://openreview.net/forum?id=0xh5CstDJU>.
- 435 [22] Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum
436 entropy exploration. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of
437 the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine
438 Learning Research*, pages 2681–2691. PMLR, 09–15 Jun 2019. URL [https://proceedings.
439 mlr.press/v97/hazan19a.html](https://proceedings.mlr.press/v97/hazan19a.html).
- 440 [23] Mikael Henaff, Roberta Raileanu, Minqi Jiang, and Tim Rocktäschel. Exploration via elliptical
441 episodic bonuses. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun
442 Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=Xg-yZos9qJQ>.
- 444 [24] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with
445 gaussian error linear units, 2017. URL <https://openreview.net/forum?id=BkOMRI51g>.
- 446 [25] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko,
447 Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High
448 definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- 449 [26] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
450 1735–1780, 1997.
- 451 [27] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer
452 and super-resolution. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 694–711. Springer, 2016.
453

- 454 [28] Łukasz Kaiser, Mohammad Babaeizadeh, Piotr Miłoś, Blażej Osieński, Roy H Campbell, Konrad
455 Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin,
456 Ryan Sepassi, George Tucker, and Henryk Michalewski. Model based reinforcement learning
457 for atari. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S1xCPJHtDB>.
458
- 459 [29] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers
460 are RNNs: Fast autoregressive transformers with linear attention. In Hal Daumé III and Aarti
461 Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume
462 119 of *Proceedings of Machine Learning Research*, pages 5156–5165. PMLR, 13–18 Jul 2020.
463 URL <https://proceedings.mlr.press/v119/katharopoulos20a.html>.
- 464 [30] Isaac Kauvar, Chris Doyle, Linqi Zhou, and Nick Haber. Curious replay for model-based
465 adaptation. In *Proceedings of the 40th International Conference on Machine Learning*, ICML’23.
466 JMLR.org, 2023.
- 467 [31] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scal-
468 able predictive uncertainty estimation using deep ensembles. In I. Guyon, U. Von
469 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, edi-
470 tors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
471 Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf.
472
- 473 [32] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther.
474 Autoencoding beyond pixels using a learned similarity metric. In Maria Florina Balcan and
475 Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine
476 Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1558–1566, New
477 York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/larsen16.html>.
478
- 479 [33] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi.
480 UNIFIED-IO: A unified model for vision, language, and multi-modal tasks. In *The Eleventh
481 International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=E01k9048soZ>.
482
- 483 [34] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek
484 Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with
485 vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer
486 Vision and Pattern Recognition (CVPR)*, pages 26439–26455, June 2024.
- 487 [35] Michael Matthews, Michael Beukman, Benjamin Ellis, Mikayel Samvelyan, Matthew Jackson,
488 Samuel Coward, and Jakob Foerster. Craftax: A lightning-fast benchmark for open-ended
489 reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2024.
- 490 [36] Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Dis-
491 covering and achieving goals via world models. *Advances in Neural Information Processing
492 Systems*, 34:24379–24391, 2021.
- 493 [37] Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world
494 models. In *The Eleventh International Conference on Learning Representations, ICLR 2023,
495 Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=vhFu1Acboxb>.
496
- 497 [38] Vincent Micheli, Eloi Alonso, and François Fleuret. Efficient world models with context-aware
498 tokenization. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna,
499 Austria, July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=BiWIERWBFX>.
500
- 501 [39] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
502 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
503 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.

- 504 [40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
505 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
506 Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
507 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
508 high-performance deep learning library. In *Advances in Neural Information Processing Systems*,
509 volume 32, 2019.
- 510 [41] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven ex-
511 ploration by self-supervised prediction. In Doina Precup and Yee Whye Teh, editors, *Pro-
512 ceedings of the 34th International Conference on Machine Learning*, volume 70 of *Pro-
513 ceedings of Machine Learning Research*, pages 2778–2787. PMLR, 06–11 Aug 2017. URL
514 <https://proceedings.mlr.press/v70/pathak17a.html>.
- 515 [42] Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions, 2018.
516 URL <https://openreview.net/forum?id=SkBYyZRZ>.
- 517 [43] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov,
518 Gabriel Barth-maroon, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom
519 Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell,
520 Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. A generalist agent. *Transactions on
521 Machine Learning Research*, 2022. ISSN 2835-8856. URL [https://openreview.net/
522 forum?id=1ikK0kHjvj](https://openreview.net/forum?id=1ikK0kHjvj). Featured Certification, Outstanding Certification.
- 523 [44] Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world
524 models are happy with 100k interactions. *arXiv preprint arXiv:2303.07109*, 2023.
- 525 [45] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
526 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF
527 Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June
528 2022.
- 529 [46] Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990-2010).
530 *IEEE Transactions on Autonomous Mental Development*, 2, 2010. ISSN 19430604. doi:
531 10.1109/TAMD.2010.2056368.
- 532 [47] Ingmar Schubert, Jingwei Zhang, Jake Bruce, Sarah Bechtle, Emilio Parisotto, Martin Ried-
533 miller, Jost Tobias Springenberg, Arunkumar Byravan, Leonard Hasenclever, and Nicolas Heess.
534 A generalist dynamics model for control. *arXiv preprint arXiv:2305.10912*, 2023.
- 535 [48] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
536 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 537 [49] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak
538 Pathak. Planning to explore via self-supervised world models. In Hal Daumé III and Aarti
539 Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume
540 119 of *Proceedings of Machine Learning Research*, pages 8583–8592. PMLR, 13–18 Jul 2020.
541 URL <https://proceedings.mlr.press/v119/sekar20a.html>.
- 542 [50] Pranav Shyam, Wojciech Jaśkowski, and Faustino Gomez. Model-based active exploration. In
543 Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International
544 Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*,
545 pages 5779–5788. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.press/v97/
546 shyam19a.html](https://proceedings.mlr.press/v97/shyam19a.html).
- 547 [51] Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang,
548 and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv
549 preprint arXiv:2307.08621*, 2023.
- 550 [52] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradi-
551 ent methods for reinforcement learning with function approximation. In S. Solla, T. Leen,
552 and K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12.
553 MIT Press, 1999. URL [https://proceedings.neurips.cc/paper_files/paper/1999/
554 file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf).

- 555 [53] Saran Tunyasuvunakool, Alistair Muldal, Yotam Doron, Siqi Liu, Steven Bohez, Josh Merel,
556 Tom Erez, Timothy Lillicrap, Nicolas Heess, and Yuval Tassa. dm_control: Software and tasks
557 for continuous control. *Software Impacts*, 6:100022, 2020. ISSN 2665-9638. doi: <https://doi.org/10.1016/j.simpa.2020.100022>. URL <https://www.sciencedirect.com/science/article/pii/S2665963820300099>.
559
- 560 [54] Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are
561 real-time game engines. *CoRR*, abs/2408.14837, 2024. URL <https://doi.org/10.48550/arXiv.2408.14837>.
562
- 563 [55] Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural discrete representation
564 learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and
565 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran
566 Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/7a98af17e63a0ac09ce2e96d03992fbc-Paper.pdf.
567
- 568 [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
569 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,
570 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, ed-
571 itors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
572 Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
573
- 574 [57] Shengjie Wang, Shaohuai Liu, Weirui Ye, Jiacheng You, and Yang Gao. Efficientzero v2: Mas-
575 tering discrete and continuous control with limited data. In *Forty-first International Conference*
576 *on Machine Learning*, 2024. URL <https://openreview.net/forum?id=LHGMXcr6zx>.
- 577 [58] ZiRui Wang, Yue DENG, Junfeng Long, and Yin Zhang. Parallelizing model-based rein-
578 forcement learning over the sequence length. In *The Thirty-eighth Annual Conference on*
579 *Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=R6N9AGyz13>.
580
- 581 [59] Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye HAO, and Mingsheng Long.
582 ivideoGPT: Interactive videoGPTs are scalable world models. In *The Thirty-eighth Annual*
583 *Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=4TENzBftZR>.
584
- 585 [60] Weipu Zhang, Gang Wang, Jian Sun, Yetian Yuan, and Gao Huang. Storm: Efficient stochastic
586 transformer based world models for reinforcement learning. *Advances in Neural Information*
587 *Processing Systems*, 36, 2024.

588 **A Models and Hyperparameters**

589 **A.1 Hyperparameters**

590 We detail shared hyperparameters in Table 1, training hyperparameters in Table 2, world model
 591 hyperparameters in Table 3, and controller hyperparameters in Table 4. Environment hyperparameters
 592 are detailed in Table 5 (Atari-100K) and Table 6 (DMC).

593 For the DMC benchmark, we use a lower embedding dimension (Table 3) due to the significantly
 594 lower dimensionality of its observations compared to other benchmarks. As in prior work (e.g.,
 595 DreamerV3), we adopt a smaller model for this setting. Note that the reduced number of Retention
 596 heads ensures a consistent head dimensionality (64).

597 Additionally, we used a limited decay range in DMC (Table 3) as observations effectively represent
 598 full MDP states, eliminating the need for long-term memory. By constraining the decay range, we
 599 explicitly encode this inductive bias into the model.

600 In the Craftax benchmark, we reduce the number of layers (Table 3) to lower computational cost. The
 601 interaction budget in Craftax is 1M steps, resulting in a particularly expensive training process. Here,
 602 we increase the decay range as an inductive bias to encourage long-term memory.

603 The weighting of intrinsic versus extrinsic rewards (Table 4) varies across benchmarks due to
 604 differences in reward structure and scale. For instance, DMC provides dense rewards with typical
 605 task scores reaching around 1000, whereas Craftax has extremely sparse rewards, with cumulative
 606 scores rarely exceeding 20—even over thousands of steps.

607 **Tuning** All remaining hyperparameters were tuned empirically or based on REM [11], with minimal
 608 impact on training cost and adjusted primarily for performance. Due to limited computational
 609 resources, we were unable to conduct extensive tuning, and we believe that further optimization could
 610 improve Simulus’s performance.

Table 1: Shared hyperparameters.

Description	Symbol	Value
Eval sampling temperature		0.5
Optimizer		AdamW
Learning rate ($\mathcal{V}, \mathcal{M}, \mathcal{C}$)		(1e-4, 2e-4, 2e-4)
AdamW β_1		0.9
AdamW β_2		0.999
Gradient clipping threshold ($\mathcal{V}, \mathcal{M}, \mathcal{C}$)		(10, 3, 3)
Weight decay ($\mathcal{V}, \mathcal{M}, \mathcal{C}$)		(0.01, 0.05, 0.01)
Prioritized replay fraction	α	0.3
Prioritized replay initial loss value	ν_0	10
\mathcal{M} ensemble size	N_{ens}	4
HL-Gauss num bins		128
Label smoothing	σ	$\frac{3}{4}\text{bin_width} = 0.1758$

Table 2: Training hyperparameters.

Description	Symbol	Atari-100K	DMC	Craftax
Horizon	H	10	20	20
Observation sequence length	K	64	3-24	147
Action sequence length	K_a	1	1-6	1
Tokenizer vocabulary size	N	512	125	(37,5,40,20,4,125)
Epochs		600	1000	10000
Experience collection epochs		500	1000	10000
Environment steps per epoch		200	500	100
Batch size ($\mathcal{V}, \mathcal{M}, \mathcal{C}$)		(128, 32, 128)	(-, 16, 128)	(-, 8, 128)
Training steps per epoch ($\mathcal{V}, \mathcal{M}, \mathcal{C}$)		(200, 200, 80)	(-, 300, 100)	(-, 100, 50)
Training start after epoch ($\mathcal{V}, \mathcal{M}, \mathcal{C}$)		(5, 25, 50)	(-, 15, 20)	(-, 250, 300)

Table 3: World model (\mathcal{M}) hyperparameters.

Description	Symbol	Atari-100K	DMC	Craftax
Number of layers		10	10	5
Number of heads		4	3	4
Embedding dimension	d	256	192	256
Dropout		0.1	0.1	0.1
Retention decay range	$[\eta_{\min}, \eta_{\max}]$	[4, 16]	[2, 2]	[8, 40]

Table 4: Actor-critic (\mathcal{C}) hyperparameters.

Description	Symbol	Atari-100K	DMC	Craftax
Environment reward weight	w^{ext}	1	1	100
Intrinsic reward weight	w^{int}	1	10	1
Encoder MLP (g_ψ) hidden layer sizes		[512]	[384]	[512, 512]
Shared backbone		True	False	True
Number of quantization values (continuous actions)			51	
(2D) Categoricals embedding dimension				64

Table 5: Atari 100K hyperparameters.

Description	Symbol	Value
Frame resolution		64×64
Frame Skip		4
Max no-ops (train, test)		(30, 1)
Max episode steps (train, test)		(20K, 108K)
Terminate on live loss (train, test)		(No, Yes)

Table 6: DeepMind Control Suite Proprioception hyperparameters.

Description	Symbol	Value
Action repeat		2

611 **A.2 The Representation Module \mathcal{V}**

612 **A.2.1 Image Observations**

613 Image observations are tokenized using a vector-quantized variational auto-encoder (VQ-VAE)
614 [55, 15]. A VQ-VAE comprises a convolutional neural network (CNN) encoder, an embedding table
615 $\mathbf{E} \in \mathbb{R}^{n \times d}$, and a CNN decoder. Here, the size of the embedding table n determines the vocabulary
616 size.

617 The encoder’s output $\mathbf{h} \in \mathbb{R}^{W \times H \times d}$ is a grid of $W \times H$ multi-channel vectors of dimension d that
618 encode high-level learned features. Each such vector is mapped to a discrete token by finding the
619 closest embedding in \mathbf{E} :

$$z = \arg \min_i \|\mathbf{h} - \mathbf{E}(i)\|,$$

620 where $\mathbf{E}(i)$ is the i -th row of \mathbf{E} . To reconstruct the original image, the decoder first maps \mathbf{z} to their
621 embeddings using \mathbf{E} . During training, the straight-through estimator [6] is used for backpropagating
622 the learning signal from the decoder to the encoder: $\hat{\mathbf{h}} = \mathbf{h} + \text{sg}(\mathbf{E}_z - \mathbf{h})$. The architecture of the
623 encoder and decoder models is presented in Table 7.

624 The optimization objective is given by

$$\mathcal{L}(\text{enc}, \text{dec}, \mathbf{E}) = \|\mathbf{o} - \text{dec}(z)\|_2^2 + \|\text{sg}(\text{enc}(\mathbf{o}) - \mathbf{E}(z))\|_2^2 + \|\text{sg}(\mathbf{E}(z) - \text{enc}(\mathbf{o}))\|_2^2 + \mathcal{L}_{\text{perceptual}}(\mathbf{o}, \text{dec}(z)),$$

625 where $\mathcal{L}_{\text{perceptual}}$ is a perceptual loss [27, 32], proposed in [37].

626 Crucially, the learned embedding table \mathbf{E} is used for embedding the (image) tokens across all stages
627 of the algorithm.

628 **A.2.2 Continuous Vectors**

629 The quantization of each feature uses 125 values (vocabulary size) in the range $[-6, 6]$, where 63
630 values are uniformly distributed in $[-\ln(1 + \pi), \ln(1 + \pi)]$ and the rest are uniformly distributed in
631 the remaining intervals.

Table 7: The encoder and decoder architectures of the VQ-VAE model. “Conv(a,b,c)” represents a convolutional layer with kernel size $a \times a$, stride of b and padding c . A value of $c = \text{Asym.}$ represents an asymmetric padding where a padding of 1 is added only to the right and bottom ends of the image tensor. “GN” represents a GroupNorm operator with 8 groups, $\epsilon = 1e - 6$ and learnable per-channel affine parameters. SiLU is the Sigmoid Linear Unit activation [24, 42]. “Interpolate” uses PyTorch’s interpolate method with scale factor of 2 and the “nearest-exact” mode.

Module	Output Shape
Encoder	
Input	$3 \times 64 \times 64$
Conv(3, 1, 1)	$64 \times 64 \times 64$
EncoderBlock1	$128 \times 32 \times 32$
EncoderBlock2	$256 \times 16 \times 16$
EncoderBlock3	$512 \times 8 \times 8$
GN	$512 \times 8 \times 8$
SiLU	$512 \times 8 \times 8$
Conv(3, 1, 1)	$256 \times 8 \times 8$
EncoderBlock	
Input	$c \times h \times w$
GN	$c \times h \times w$
SiLU	$c \times h \times w$
Conv(3, 2, Asym.)	$2c \times \frac{h}{2} \times \frac{w}{2}$
Decoder	
Input	$256 \times 8 \times 8$
BatchNorm	$256 \times 8 \times 8$
Conv(3, 1, 1)	$256 \times 8 \times 8$
DecoderBlock1	$128 \times 16 \times 16$
DecoderBlock2	$64 \times 32 \times 32$
DecoderBlock3	$64 \times 64 \times 64$
GN	$64 \times 64 \times 64$
SiLU	$64 \times 64 \times 64$
Conv(3, 1, 1)	$3 \times 64 \times 64$
DecoderBlock	
Input	$c \times h \times w$
GN	$c \times h \times w$
SiLU	$c \times h \times w$
Interpolate	$c \times 2h \times 2w$
Conv(3, 1, 1)	$\frac{c}{2} \times 2h \times 2w$

632 **A.3 The World Model \mathcal{M}**

633 **Embedding Details** Each token in $\mathbf{z}^{(i)}$ of each modality is mapped to a d -dimensional embedding
 634 vector $\mathbf{X}^{(i)}$ using the embedding (look-up) table $\mathbf{E}^{(i)}$ of modality κ_i . The embedding vector that
 635 corresponds to token z is simply the z -th row in the embedding table. Formally, $\mathbf{x}_{t,j}^{(i)} = \mathbf{E}^{(i)}(l)$, $l =$
 636 $z_{t,j}^{(i)}$ where $\mathbf{E}(l)$ refers to the l -th row in \mathbf{E} . In the special case of 2D categorical inputs, $\mathbf{x}_{t,j}^{(i)} =$
 637 $\frac{1}{C} \sum_{n=1}^C \mathbf{E}_n^{(i)}(l_n)$, $l_n = z_{t,j,n}^{(i)}$ where C is the number of channels and i is the index of the 2D
 638 categorical modality in κ .

639 To concatenate the embeddings, we use the following order among the modalities: images, continuous
 640 vectors, categorical variables, and 2D categoricals.

641 **Prediction Heads** Each prediction head in \mathcal{M} is a multi-layer perceptron (MLP) with a single
 642 hidden layer of dimension $2d$ where d is the embedding dimension.

643 **Epistemic Uncertainty Estimation** Working with discrete distributions enables efficient
 644 entropy computation and ensures that the ensemble disagreement term δ_t is bounded by
 645 $\frac{1}{|\mathbf{z}_t|} \sum_{z \in \mathbf{z}_t} \log(\text{vocab_size}(z))$.

646 **A.3.1 Optimization**

647 For each training example in the form of a trajectory segment in token representation $\tau =$
 648 $\mathbf{z}_1, \mathbf{z}_1^a, \dots, \mathbf{z}_H, \mathbf{z}_H^a$, the optimization objective is given by

$$\mathcal{L}_{\mathcal{M}}(\theta, \phi, \tau) = \sum_{t=1}^H \mathcal{L}_{\text{obs}}(\theta, \mathbf{z}_t, p_{\theta}(\hat{\mathbf{z}}_t | \mathbf{Y}_t^u)) + \mathcal{L}_{\text{reward}}(\theta, r_t, \hat{r}_t) - \log(p_{\theta}(d_t | \mathbf{Y}_t^u)) \\ + \sum_{i=1}^{N_{\text{ens}}} \mathcal{L}_{\text{obs}}(\phi_i, \mathbf{z}_t, p_{\phi_i}(\hat{\mathbf{z}}_t | \text{sg}(\mathbf{Y}_t^u))),$$

649 where

$$\mathcal{L}_{\text{obs}}(\theta, \mathbf{z}_t, p_{\theta}(\hat{\mathbf{z}}_t | \mathbf{Y}_t^u)) = -\frac{1}{K} \sum_{i=1}^K \log p_{\theta}(z_i | \mathbf{y}_i)$$

650 is the average of the cross-entropy losses of the individual tokens, and $\mathcal{L}_{\text{reward}}(\theta, r_t, \hat{r}_t)$ is the $\mathcal{L}_{\text{HL-Gauss}}$
 651 loss with the respective parameters of the reward head. Here, \mathbf{y}_i is the vector of \mathbf{Y}_t^u that corresponds
 652 to z_i , the i -th token of \mathbf{z}_t .

653 **A.3.2 Retentive Networks**

654 Retentive Networks (RetNet) [51] are sequence model architectures with a Transformer-like structure
 655 [56]. However, RetNet replaces the self-attention mechanism with a linear-attention [29] based
 656 Retention mechanism. At a high level, given an input sequence $\mathbf{X} \in \mathbb{R}^{|\mathbf{X}| \times d}$ of d -dimensional
 657 vectors, the Retention operator outputs

$$\text{Retention}(\mathbf{X}) = (\mathbf{Q}\mathbf{K}^{\top} \odot \mathbf{D})\mathbf{V},$$

658 where $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are the queries, keys, and values, respectively, and \mathbf{D} is a causal mask and decay
 659 matrix. Notably, the softmax operation is discarded in Retention and other linear attention methods.
 660 As a linear attention method, the computation can also be carried in a recurrent form:

$$\text{Retention}(\mathbf{x}_t, \mathbf{S}_{t-1}) = \mathbf{S}_t \mathbf{q}_t, \\ \mathbf{S}_t = \eta \mathbf{S}_{t-1} + \mathbf{v}_t \mathbf{k}_t^{\top} \in \mathbb{R}^{d \times d},$$

661 where η is a decay factor, \mathbf{S}_t is a recurrent state, and $\mathbf{S}_0 = 0$. In addition, a hybrid form of recurrent
 662 and parallel forward computation known as the chunkwise mode allows to balance the quadratic
 663 cost of the parallel form and the sequential cost of the recurrent form by processing the input as a
 664 sequence of chunks. We refer the reader to [51] for the full details about this architecture.

665 In our implementation, since inputs are complete observation-action block sequences $\mathbf{X}_1, \dots, \mathbf{X}_t$,
 666 we configure the decay factors of the multi-scale retention operator in block units:

$$\eta = 1 - 2^{-\text{linspace}(\log_2(K\eta_{\min}), \log_2(K\eta_{\max}), N_h)},$$

667 where $\text{linspace}(a, b, c)$ is a sequence of c values evenly distributed between a and b , N_h is the
 668 number of retention heads, and η_{\min}, η_{\max} are hyperparameters that control the memory decay in
 669 observation-action block units.

670 A.3.3 Parallel Observation Prediction (POP)

671 POP [11] is a mechanism for parallel generation of non-causal subsequences such as observations in
 672 token representation. It’s purpose is to improve generation efficiency by alleviating the sequential
 673 bottleneck caused by generating observations a single token at a time (as done in language models).
 674 However, to achieve this goal, POP also includes a mechanism for maintaining training efficiency.
 675 Specifically, POP extends the chunkwise forward mode of RetNet to maintain efficient training of the
 676 sequence model.

677 To generate multiple tokens into the future at once, POP introduces a set of prediction tokens
 678 $\mathbf{u} = u_1, \dots, u_K$ and embeddings $\mathbf{X}^u \in \mathbb{R}^{K \times d}$ where K is the number of tokens in an observation.
 679 Each token in \mathbf{u} corresponds to an observation token in \mathbf{z} . These tokens, and their respective learned
 680 embeddings, serve as a learned prior.

681 Let $\mathbf{X}_1, \dots, \mathbf{X}_T$ be a sequence of T observation-action (embeddings) blocks. Given \mathbf{S}_{t-1} sum-
 682 marizing all key-value outer products of elements of $\mathbf{X}_{\leq t-1}$, the outputs \mathbf{Y}^u from which the next
 683 observation tokens are predicted are given by:

$$(\cdot, \mathbf{Y}_t^u) = f_\theta(\mathbf{S}_{t-1}, \mathbf{X}^u).$$

684 Importantly, the recurrent state is never updated based on the prediction tokens \mathbf{u} (or their embeddings).
 685 The next observation tokens $\hat{\mathbf{z}}_t$ are sampled from $p_\theta(\hat{\mathbf{z}}_t | \mathbf{Y}_t^u)$. Then, the next action is generated by the
 686 controller, and the next observation-action block \mathbf{X}_t can be processed to predict the next observation
 687 $\hat{\mathbf{z}}_{t+1}$.

688 To maintain efficient training, a two step computation is carried at each RetNet layer. First, all
 689 recurrent states \mathbf{S}_t for all $1 \leq t \leq T$ are calculated in parallel. Although there is an auto-regressive
 690 relationship between time steps, the linear structure of \mathbf{S} allows to calculate the compute-intensive
 691 part of each state in parallel and incorporate past information efficiently afterwards. In the second
 692 step, all outputs \mathbf{Y}_t^u for all $1 \leq t \leq T$ are computed in parallel, using the appropriate states \mathbf{S}_{t-1}
 693 and \mathbf{X}^u in batch computation. Note that this computation involves delicate positional information
 694 handling. We refer the reader to [11] for full details of this computation.

695 **A.4 The Controller \mathcal{C}**

696 **Critic** The value prediction uses 128 bins in the range $\mathbf{b} = (-15, \dots, 15)$.

697 **Continuous Action Spaces** The policy network outputs $m = 51$ logits corresponding to m
 698 quantization values uniformly distributed in $[-1, 1]$ for each individual action in the action vector.

699 **A.4.1 Input Encoding**

700 The controller \mathcal{C} operates in the latent token space. Token trajectories $\tau = \mathbf{z}_1, \mathbf{z}_1^a, \dots, \mathbf{z}_H, \mathbf{z}_H^a$ are
 701 processed sequentially by the LSTM model. At each time step t , the network gets \mathbf{z}_t as input, outputs
 702 $\pi(\mathbf{a}_t | \tau_{\leq t-1}, \mathbf{z}_t)$ and $\hat{V}^\pi(\mathbf{a}_t | \tau_{\leq t-1}, \mathbf{z}_t)$, samples an action \mathbf{a}_t and then process the sampled action as
 703 another sequence element.

704 The processing of actions involve embedding the action into a latent vector which is then provided
 705 as input to the LSTM. Embedding of continuous action tokens is performed by first reconstructing
 706 the continuous action vector and then computing the embedding using a linear projection. Discrete
 707 tokens are embedded using a dedicated embedding table.

708 To embed observation tokens \mathbf{z} , the tokens of each modality are processed by a modality-specific
 709 encoder. The outputs of the encoders are concatenated and further processed by a MLP g_ψ that
 710 combines the information into a single vector latent representation.

711 The image encoder is a convolutional neural network (CNN). Its architecture is given in Table 8.

712 Categorical variables are embedded using a learned embedding table. For 2D categoricals, shared
 713 per-channel embedding tables map tokens to embedding vectors, which are averaged to obtain a
 714 single embedding for each multi-channel token vector. For both types of categorical inputs we use 64
 715 dimensional embeddings. The embeddings are concatenated and processed by g_ψ .

Table 8: The image observation encoder architecture of the actor-critic controller \mathcal{C} .

Module	Output Shape
Input	$256 \times 8 \times 8$
Conv(3, 1, 1)	$128 \times 8 \times 8$
SiLU	$128 \times 8 \times 8$
Conv(3, 1, 1)	$64 \times 8 \times 8$
SiLU	$64 \times 8 \times 8$
Flatten	4096
Linear	512
SiLU	512

716 **A.4.2 Optimization**

717 λ -returns are computed for each generated trajectory segment $\hat{\tau} =$
 718 $(\mathbf{z}_1, \mathbf{a}_1, \bar{r}_1, d_1, \hat{\mathbf{z}}_2, \mathbf{a}_2, \bar{r}_2, d_2, \dots, \hat{\mathbf{z}}_H, \mathbf{a}_H, \bar{r}_H, d_H)$:

$$G_t = \begin{cases} \bar{r}_t + \gamma(1 - d_t)((1 - \lambda)\hat{V}_{t+1}^\pi + \lambda G_{t+1}) & t < H \\ \hat{V}_H^\pi & t = H \end{cases}$$

719 where $\hat{V}_t^\pi = \hat{V}^\pi(\hat{\tau}_{\leq t})$. These λ -returns are used as targets for critic learning. For policy learning, a
 720 REINFORCE [52] objective is used, with a normalized \hat{V}^π baseline for variance reduction:

$$\mathcal{L}_\pi(\psi) = \mathbb{E}_\pi \left[\sum_{t=1}^H \text{sg} \left(\frac{G_t - \hat{V}_t^\pi}{\max(1, c)} \right) \log \pi(\mathbf{a}_t | \hat{\tau}_{\leq t-1}, \hat{\mathbf{z}}_t) + w_{\text{ent}} \mathcal{H}(\pi(\mathbf{a}_t | \hat{\tau}_{\leq t-1}, \hat{\mathbf{z}}_t)) \right],$$

721 where c is an estimate of the effective return scale similar to DreamerV3 [20] and w_{ent} is a hyper-
 722 parameter that controls the entropy regularization weight. c is calculated as the difference between
 723 the running average estimators of the 97.5 and 2.5 return percentiles, based on a window of return
 724 estimates obtained in the last 500 batches (imagination).

725 **B Additional Results**

726 The average per-game scores for Atari-100K are presented in Table 9. The performance profile plot
 727 for Atari 100K is presented in Figure 7.

Table 9: Mean returns on the 26 games of the Atari 100k benchmark followed by averaged human-normalized performance metrics. Each game score is computed as the average of 5 runs with different seeds. Bold face mark the best score.

Game	Random	Human	DreamerV3	TWM	STORM	DIAMOND	REM	SIMULUS (ours)
Alien	227.8	7127.7	959.4	674.6	983.6	744.1	607.2	687.2
Amidar	5.8	1719.5	139.1	121.8	204.8	225.8	95.3	102.4
Assault	222.4	742.0	705.6	682.6	801.0	1526.4	1764.2	1822.8
Asterix	210.0	8503.3	932.5	1116.6	1028.0	3698.5	1637.5	1369.1
BankHeist	14.2	753.1	648.7	466.7	641.2	19.7	19.2	347.1
BattleZone	2360.0	37187.5	12250.0	5068.0	13540.0	4702.0	11826.0	13262.0
Boxing	0.1	12.1	78.0	77.5	79.7	86.9	87.5	93.5
Breakout	1.7	30.5	31.1	20.0	15.9	132.5	90.7	148.9
ChopperCommand	811.0	7387.8	410.0	1697.4	1888.0	1369.8	2561.2	3611.6
CrazyClimber	10780.5	35829.4	97190.0	71820.4	66776.0	99167.8	76547.6	93433.2
DemonAttack	152.1	1971.0	303.3	350.2	164.6	288.1	5738.6	4787.6
Freeway	0.0	29.6	0.0	24.3	0.0	33.3	32.3	31.9
Frostbite	65.2	4334.7	909.4	1475.6	1316.0	274.1	240.5	258.4
Gopher	257.6	2412.5	3730.0	1674.8	8239.6	5897.9	5452.4	4363.2
Hero	1027.0	30826.4	11160.5	7254.0	11044.3	5621.8	6484.8	7466.8
Jamesbond	29.0	302.8	444.6	362.4	509.0	427.4	391.2	678.0
Kangaroo	52.0	3035.0	4098.3	1240.0	4208.0	5382.2	467.6	6656.0
Krull	1598.0	2665.5	7781.5	6349.2	8412.6	8610.1	4017.7	6677.3
KungFuMaster	258.5	22736.3	21420.0	24554.6	26182.0	18713.6	25172.2	31705.4
MsPacman	307.3	6951.6	1326.9	1588.4	2673.5	1958.2	962.5	1282.7
Pong	-20.7	14.6	18.4	18.8	11.3	20.4	18.0	19.9
PrivateEye	24.9	69571.3	881.6	86.6	7781.0	114.3	99.6	100.0
Qbert	163.9	13455.0	3405.1	3330.8	4522.5	4499.3	743.0	2425.6
RoadRunner	11.5	7845.0	15565.0	9109.0	17564.0	20673.2	14060.2	24471.8
Seaquest	68.4	42054.7	618.0	774.4	525.2	551.2	1036.7	1800.4
UpNDown	533.4	11693.2	7567.1	15981.7	7985.0	3856.3	3757.6	10416.5
#Superhuman (\uparrow)	0	N/A	9	8	9	11	12	13
Mean (\uparrow)	0.000	1.000	1.124	0.956	1.222	1.459	1.222	1.645
Median (\uparrow)	0.000	1.000	0.485	0.505	0.425	0.373	0.280	0.982
IQM (\uparrow)	0.000	1.000	0.487	0.459	0.561	0.641	0.673	0.990
Optimality Gap (\downarrow)	1.000	0.000	0.510	0.513	0.472	0.480	0.482	0.412

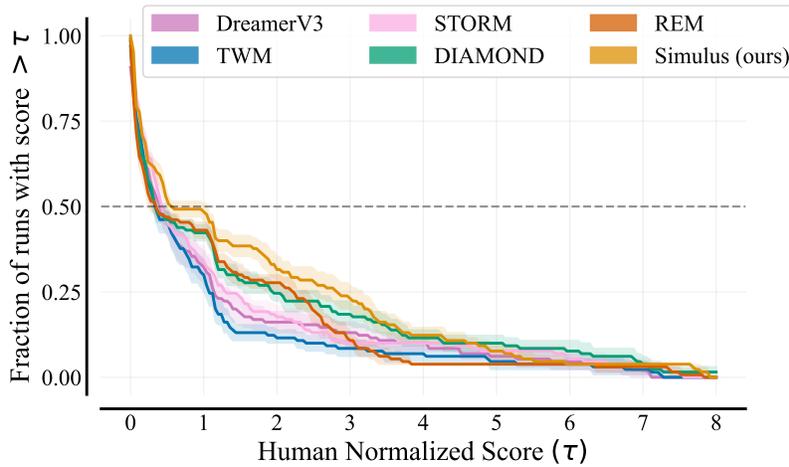


Figure 7: Performance profile. For each human-normalized score value on the x-axis, the curve of each algorithm represents the fraction of runs achieving a score greater than that value. Shaded regions denote pointwise 95% confidence intervals, computed using stratified bootstrap sampling [2].

728 **B.1 Interfering Objectives in RSSM Optimization**

729 Here, we study the interplay between the objectives of a Dreamer-like world model, PWM [58], which
 730 uses a slightly modified version of the recurrent state space model (RSSM) of Dreamer. Concretely,
 731 we aim to understand whether the representation and sequence modeling objectives interfere by
 732 decoupling the optimization of the encoder-decoder models from that of the world model. We opted
 733 for this implementation due to its simplicity, fast runtime, and accessibility, as it is written in PyTorch
 734 [40].

735 Formally, the model consists of the following components:

$$\begin{aligned} \text{Encoder: } z_t &\sim q_\theta(z_t|o_t), \\ \text{Decoder: } \hat{o}_t &\sim p_\theta(\hat{o}_t|z_t), \\ \text{Sequence model: } h_t, x_t &= f_\theta(x_{t-1}, z_{t-1}, a_{t-1}), \\ \text{Dynamics predictor: } \hat{z}_t &\sim p_\theta(\hat{z}_t|h_t). \end{aligned}$$

736 We omit the reward and termination predictors and objectives for brevity. Note that in contrast to the
 737 RSSM in Dreamer, the encoder and decoder models do not depend on the recurrent state h_t, x_t . The
 738 optimization objective of PWM is given by

$$\mathcal{L}(\theta) = \mathbb{E}_{q_\theta} \left[\sum_{t=1}^T \beta_{\text{pred}} \mathcal{L}_{\text{pred}}(\theta) + \beta_{\text{dyn}} \mathcal{L}_{\text{dyn}}(\theta) + \beta_{\text{rep}} \mathcal{L}_{\text{rep}}(\theta) \right],$$

739 where $\beta_{\text{pred}}, \beta_{\text{dyn}}$, and β_{rep} are coefficients and

$$\begin{aligned} \mathcal{L}_{\text{pred}}(\theta) &= \|\hat{o}_t - o_t\|_2^2, \\ \mathcal{L}_{\text{dyn}}(\theta) &= \max(1, \text{KL}[\text{sg}(q_\theta(z_t|o_t)) \| p_\theta(\hat{z}_t|h_t)]), \\ \mathcal{L}_{\text{rep}}(\theta) &= \max(1, \text{KL}[q_\theta(z_t|o_t) \| \text{sg}(p_\theta(\hat{z}_t|h_t))]). \end{aligned}$$

740 To decouple the optimization, we modify the sequence model by introducing a stop-gradient operator
 741 on the encoder’s output during world model training:

$$h_t, x_t = f_\theta(x_{t-1}, \text{sg}(z_{t-1}), a_{t-1}).$$

742 Moreover, this modification allows to train the encoder-decoder models using large batches of single
 743 frames, rather than small highly-correlated batches of long trajectories. This further highlights the
 744 flexibility and advantage of a modular design.

745 We compare the original PWM algorithm to its decoupled variant, PWM-decoupled, across four Atari
 746 environments: Breakout, DemonAttack, Hero, and RoadRunner. These are games where PWM
 747 performed either particularly well (Hero and RoadRunner) or poorly (Breakout and DemonAttack).
 748 Each variant is trained online from scratch on each game. The results are presented in Figure 8.
 749 In addition, we present the achieved episodic returns in Figure 9, and the reconstruction quality of
 750 example episodes in Figure 10 (PWM) and Figure 11 (PWM-decoupled).

751 Although our results are based on a single random seed and are limited to only four environments, we
 752 observe a consistent trend. First, the reconstruction losses are consistently and significantly lower
 753 when decoupling the optimization, while the dynamics losses are significantly higher. This suggests
 754 that the objectives are interfering.

755 Second, we observe similar or better episodic returns (Figure 9) using the decoupled optimization,
 756 suggesting that the higher dynamics loss might not lead to worse world modeling performance in
 757 practice. Note that a higher dynamics loss in this case does not necessarily mean worse performance,
 758 as for example multiple discrete combinations could represent the same or similar frame. Thus, when
 759 the dynamics model fails to predict a specific combination, it leads to high loss values while the
 760 underlying representations are accurate.

761 Lastly, we report that similar trends were observed when training only the world model in an
 762 offline, supervised-learning fashion on pre-collected datasets. We explored this setting to eliminate
 763 complexities that may arise due to the online collection of the data.

764 While the presented preliminary results are noisy and limited, we believe that they uncover an
 765 interesting observation on the design and optimization of current world models.

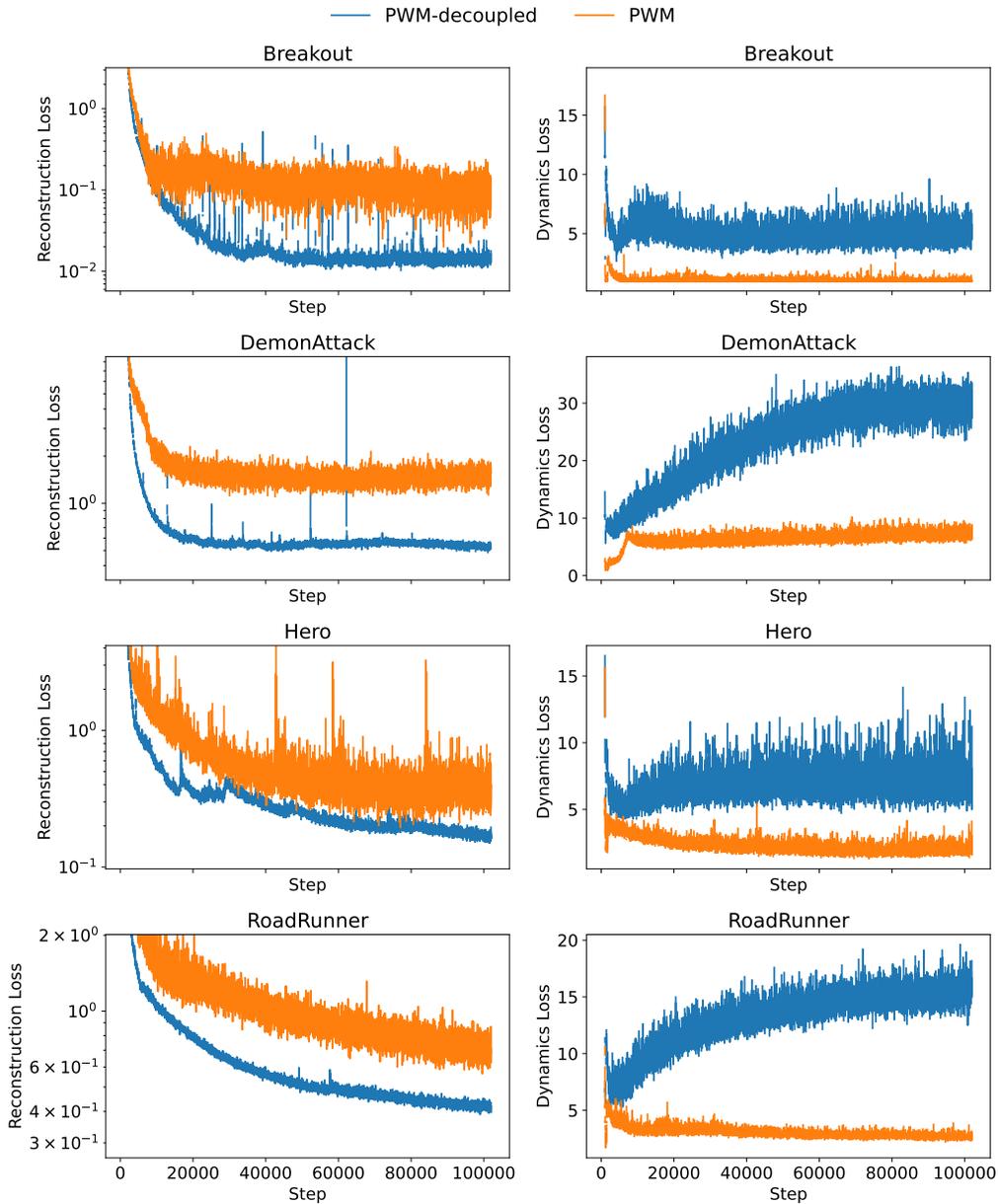


Figure 8: Reconstruction ($\mathcal{L}_{\text{pred}}$) and dynamics (\mathcal{L}_{dyn}) losses of PWM and PWM-decoupled on four Atari games (single seed). The first column uses a log-scaled y-axis. Decoupling the optimization objectives consistently reduces reconstruction loss while increasing dynamics loss, suggesting interference between the two objectives.

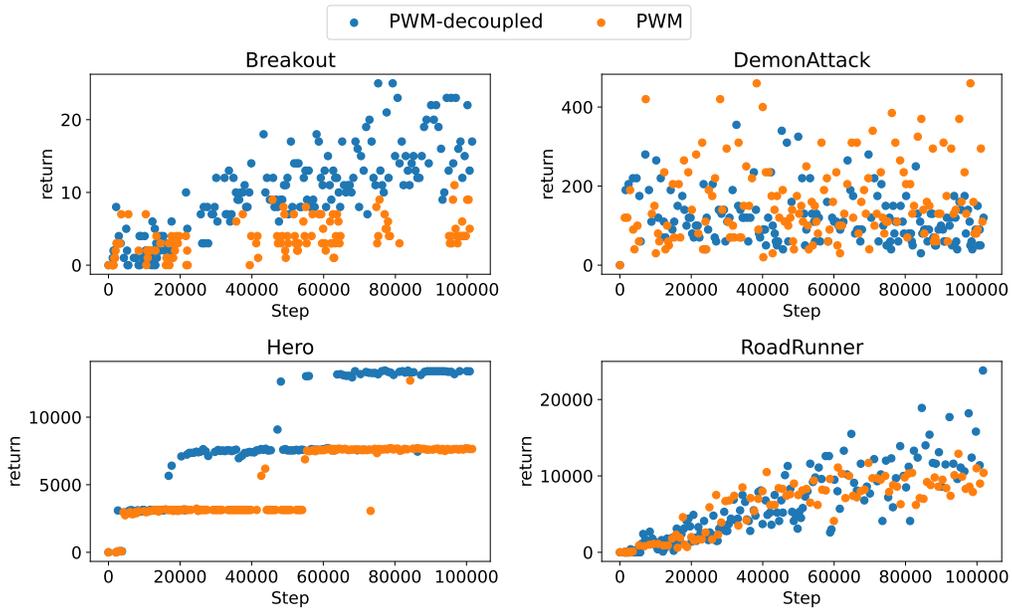


Figure 9: Agent episodic returns throughout training of PWM and PWM-decoupled on four Atari games (single seed). Each marker corresponds to a single episode.

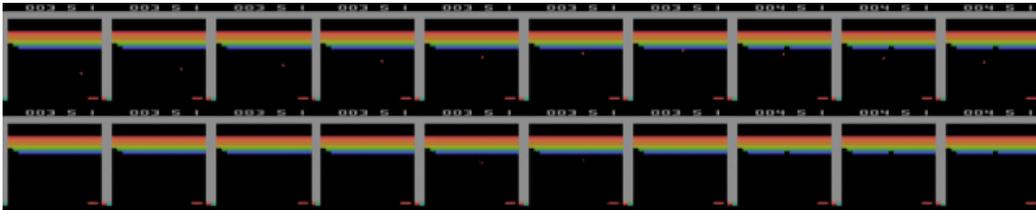


Figure 10: Ground truth (top) and reconstructed (bottom) frames from a training episode of PWM after 50K steps (half way through training). Notably, the ball is missing in most frames, suggesting the reason for its poor performance in this game.

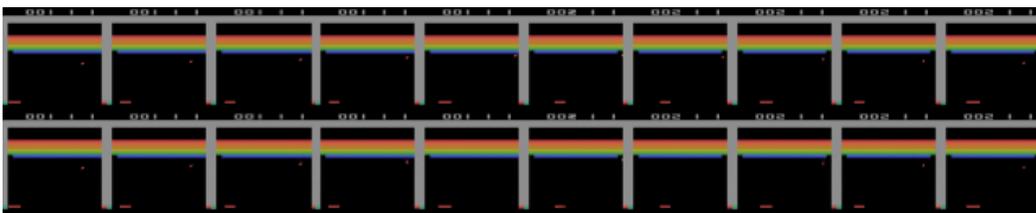


Figure 11: Ground truth (top) and reconstructed (bottom) frames from a training episode of PWM-decoupled after 50K steps (half way through training). Here, the ball is reconstructed in most frames, demonstrating the significantly improved representation performance.

766 C Implementation Details

767 **Ablation Studies** For the Atari 100K benchmark, we conducted ablations on the
768 following games: Assault, Breakout, ChopperCommand, CrazyClimber, JamesBond,
769 Kangaroo, Seaquest, and UpNDown. For the DeepMind Control Suite, we used the
770 tasks: acrobot-swingup, cartpole-swingup-sparse, cheetah-run, finger-turn-hard,
771 hopper-stand, pendulum-swingup, reacher-hard, and walker-run.

772 **Code** We open-source our code and trained model weights. Our code is written in Pytorch [40].

773 **Hardware** All Atari and DMC experiments were performed on V100 GPUs, while for Craftax a
774 single RTX 4090 was used.

775 **Run Times** Experiments on Atari require approximately 12 hours on an RTX 4090 GPU and
776 around 29 hours on a V100 GPU. For DMC, the runtime is about 40 hours on a V100 GPU. Craftax
777 runs take roughly 94 hours, equivalent to 3.9 days.

778 **Craftax** The official environment provides the categorical variables in one-hot encoding format.
779 Our implementation translates these variables to integer values which can be interpreted as tokens.

780 **Setup in Atari Freeway** For the Freeway environment, we adopted a modified sampling strategy
781 where a temperature of 0.01 is used instead of the standard value of 1, following [37, 11]. This
782 adjustment helps directing the agent toward rewarding paths. Note that alternative approaches in the
783 literature tackle the exploration challenge through different mechanisms, including epsilon-greedy
784 exploration schedules and deterministic action selection via argmax policies [37].

785 **NeurIPS Paper Checklist**

786 **1. Claims**

787 Question: Do the main claims made in the abstract and introduction accurately reflect the
788 paper's contributions and scope?

789 Answer: [\[Yes\]](#)

790 Justification: We provide extensive empirical evidence in Section 3, including ablation
791 studies, which directly relate to our contributions and claims. The scope of our paper
792 is sample-efficient, planning-free world model agents (RL), which is explicitly stated in
793 the abstract and introduction, while it is also reflected by the choice of baselines in our
794 experiments.

795 Guidelines:

- 796 • The answer NA means that the abstract and introduction do not include the claims
797 made in the paper.
- 798 • The abstract and/or introduction should clearly state the claims made, including the
799 contributions made in the paper and important assumptions and limitations. A No or
800 NA answer to this question will not be perceived well by the reviewers.
- 801 • The claims made should match theoretical and experimental results, and reflect how
802 much the results can be expected to generalize to other settings.
- 803 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
804 are not attained by the paper.

805 **2. Limitations**

806 Question: Does the paper discuss the limitations of the work performed by the authors?

807 Answer: [\[Yes\]](#)

808 Justification: Section 5 explicitly discuss the limitations of our work. Additional limitations
809 are discussed in Section 3 (e.g., the absence of ablations on Craftax due to computational
810 limitations).

811 Guidelines:

- 812 • The answer NA means that the paper has no limitation while the answer No means that
813 the paper has limitations, but those are not discussed in the paper.
- 814 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 815 • The paper should point out any strong assumptions and how robust the results are to
816 violations of these assumptions (e.g., independence assumptions, noiseless settings,
817 model well-specification, asymptotic approximations only holding locally). The authors
818 should reflect on how these assumptions might be violated in practice and what the
819 implications would be.
- 820 • The authors should reflect on the scope of the claims made, e.g., if the approach was
821 only tested on a few datasets or with a few runs. In general, empirical results often
822 depend on implicit assumptions, which should be articulated.
- 823 • The authors should reflect on the factors that influence the performance of the approach.
824 For example, a facial recognition algorithm may perform poorly when image resolution
825 is low or images are taken in low lighting. Or a speech-to-text system might not be
826 used reliably to provide closed captions for online lectures because it fails to handle
827 technical jargon.
- 828 • The authors should discuss the computational efficiency of the proposed algorithms
829 and how they scale with dataset size.
- 830 • If applicable, the authors should discuss possible limitations of their approach to
831 address problems of privacy and fairness.
- 832 • While the authors might fear that complete honesty about limitations might be used by
833 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
834 limitations that aren't acknowledged in the paper. The authors should use their best
835 judgment and recognize that individual actions in favor of transparency play an impor-
836 tant role in developing norms that preserve the integrity of the community. Reviewers
837 will be specifically instructed to not penalize honesty concerning limitations.

838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Our paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: In Section 2 and in the appendix we discuss our method in full detail, including architectures, hyperparameters, etc.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: In our abstract and appendix we provide a link to the code and trained model weights. Our code has a detailed readme file for easy usage, and we also provide Docker support, which enables an easy environment setup and enhances reproducibility on any operation system.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We specify all experimental details in Section 3 and in Appendix A and C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We utilize the rliable toolkit [2] to generate plots with appropriate error bars (Figure 4 bottom, Figure 6). Figure 5 also includes error bars. In addition, our open-sourced repository includes the full results of all runs.

Guidelines:

- The answer NA means that the paper does not include experiments.

- 943 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
944 dence intervals, or statistical significance tests, at least for the experiments that support
945 the main claims of the paper.
- 946 • The factors of variability that the error bars are capturing should be clearly stated (for
947 example, train/test split, initialization, random drawing of some parameter, or overall
948 run with given experimental conditions).
- 949 • The method for calculating the error bars should be explained (closed form formula,
950 call to a library function, bootstrap, etc.)
- 951 • The assumptions made should be given (e.g., Normally distributed errors).
- 952 • It should be clear whether the error bar is the standard deviation or the standard error
953 of the mean.
- 954 • It is OK to report 1-sigma error bars, but one should state it. The authors should
955 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
956 of Normality of errors is not verified.
- 957 • For asymmetric distributions, the authors should be careful not to show in tables or
958 figures symmetric error bars that would yield results that are out of range (e.g. negative
959 error rates).
- 960 • If error bars are reported in tables or plots, The authors should explain in the text how
961 they were calculated and reference the corresponding figures or tables in the text.

962 8. Experiments compute resources

963 Question: For each experiment, does the paper provide sufficient information on the com-
964 puter resources (type of compute workers, memory, time of execution) needed to reproduce
965 the experiments?

966 Answer: [Yes]

967 Justification: We provide this information in Appendix C.

968 Guidelines:

- 969 • The answer NA means that the paper does not include experiments.
- 970 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
971 or cloud provider, including relevant memory and storage.
- 972 • The paper should provide the amount of compute required for each of the individual
973 experimental runs as well as estimate the total compute.
- 974 • The paper should disclose whether the full research project required more compute
975 than the experiments reported in the paper (e.g., preliminary or failed experiments that
976 didn't make it into the paper).

977 9. Code of ethics

978 Question: Does the research conducted in the paper conform, in every respect, with the
979 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

980 Answer: [Yes]

981 Justification: Our work follows the NeurIPS Code of Ethics. No human subjects or partici-
982 pants were involved. We found no special concerns beyond those related to the general topic
983 of deep reinforcement learning.

984 Guidelines:

- 985 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 986 • If the authors answer No, they should explain the special circumstances that require a
987 deviation from the Code of Ethics.
- 988 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
989 eration due to laws or regulations in their jurisdiction).

990 10. Broader impacts

991 Question: Does the paper discuss both potential positive societal impacts and negative
992 societal impacts of the work performed?

993 Answer: [NA]

994 Justification: This paper presents a foundational work in the field of Machine Learning. As
995 such, there are no direct positive or negative societal impacts.

996 Guidelines:

- 997 • The answer NA means that there is no societal impact of the work performed.
- 998 • If the authors answer NA or No, they should explain why their work has no societal
999 impact or why the paper does not address societal impact.
- 1000 • Examples of negative societal impacts include potential malicious or unintended uses
1001 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1002 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1003 groups), privacy considerations, and security considerations.
- 1004 • The conference expects that many papers will be foundational research and not tied
1005 to particular applications, let alone deployments. However, if there is a direct path to
1006 any negative applications, the authors should point it out. For example, it is legitimate
1007 to point out that an improvement in the quality of generative models could be used to
1008 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1009 that a generic algorithm for optimizing neural networks could enable people to train
1010 models that generate Deepfakes faster.
- 1011 • The authors should consider possible harms that could arise when the technology is
1012 being used as intended and functioning correctly, harms that could arise when the
1013 technology is being used as intended but gives incorrect results, and harms following
1014 from (intentional or unintentional) misuse of the technology.
- 1015 • If there are negative societal impacts, the authors could also discuss possible mitigation
1016 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1017 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1018 feedback over time, improving the efficiency and accessibility of ML).

1019 11. Safeguards

1020 Question: Does the paper describe safeguards that have been put in place for responsible
1021 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1022 image generators, or scraped datasets)?

1023 Answer: [NA]

1024 Justification: Our work does not pose any additional risks beyond those of common deep
1025 reinforcement learning methods. Hence, we do not introduce additional safeguards.

1026 Guidelines:

- 1027 • The answer NA means that the paper poses no such risks.
- 1028 • Released models that have a high risk for misuse or dual-use should be released with
1029 necessary safeguards to allow for controlled use of the model, for example by requiring
1030 that users adhere to usage guidelines or restrictions to access the model or implementing
1031 safety filters.
- 1032 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1033 should describe how they avoided releasing unsafe images.
- 1034 • We recognize that providing effective safeguards is challenging, and many papers do
1035 not require this, but we encourage authors to take this into account and make a best
1036 faith effort.

1037 12. Licenses for existing assets

1038 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1039 the paper, properly credited and are the license and terms of use explicitly mentioned and
1040 properly respected?

1041 Answer: [Yes]

1042 Justification: Our paper cites all relevant assets, and our open-sourced repository includes
1043 a credits section with the relevant credits. We follow the licenses of all assets used in our
1044 work.

1045 Guidelines:

- 1046 • The answer NA means that the paper does not use existing assets.

- 1047 • The authors should cite the original paper that produced the code package or dataset.
- 1048 • The authors should state which version of the asset is used and, if possible, include a
- 1049 URL.
- 1050 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1051 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 1052 service of that source should be provided.
- 1053 • If assets are released, the license, copyright information, and terms of use in the
- 1054 package should be provided. For popular datasets, paperswithcode.com/datasets
- 1055 has curated licenses for some datasets. Their licensing guide can help determine the
- 1056 license of a dataset.
- 1057 • For existing datasets that are re-packaged, both the original license and the license of
- 1058 the derived asset (if it has changed) should be provided.
- 1059 • If this information is not available online, the authors are encouraged to reach out to
- 1060 the asset’s creators.

1061 13. New assets

1062 Question: Are new assets introduced in the paper well documented and is the documentation

1063 provided alongside the assets?

1064 Answer: [Yes]

1065 Justification: Our open-sourced repository includes all new assets and is well documented.

1066 Guidelines:

- 1067 • The answer NA means that the paper does not release new assets.
- 1068 • Researchers should communicate the details of the dataset/code/model as part of their
- 1069 submissions via structured templates. This includes details about training, license,
- 1070 limitations, etc.
- 1071 • The paper should discuss whether and how consent was obtained from people whose
- 1072 asset is used.
- 1073 • At submission time, remember to anonymize your assets (if applicable). You can either
- 1074 create an anonymized URL or include an anonymized zip file.

1075 14. Crowdsourcing and research with human subjects

1076 Question: For crowdsourcing experiments and research with human subjects, does the paper

1077 include the full text of instructions given to participants and screenshots, if applicable, as

1078 well as details about compensation (if any)?

1079 Answer: [NA]

1080 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1081 Guidelines:

- 1082 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1083 human subjects.
- 1084 • Including this information in the supplemental material is fine, but if the main contribu-
- 1085 tion of the paper involves human subjects, then as much detail as possible should be
- 1086 included in the main paper.
- 1087 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 1088 or other labor should be paid at least the minimum wage in the country of the data
- 1089 collector.

1090 15. Institutional review board (IRB) approvals or equivalent for research with human

1091 subjects

1092 Question: Does the paper describe potential risks incurred by study participants, whether

1093 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

1094 approvals (or an equivalent approval/review based on the requirements of your country or

1095 institution) were obtained?

1096 Answer: [NA]

1097 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1098

Guidelines:

1099

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

1100

1101

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

1102

1103

1104

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.

1105

1106

1107

- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

1108

1109

16. Declaration of LLM usage

1110

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

1111

1112

1113

1114

Answer: [NA]

1115

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

1116

1117

Guidelines:

1118

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.

1119

1120

- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

1121