

---

# Policy-shaped prediction: improving world modeling through interpretability

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Model-based reinforcement learning (MBRL) offers sample-efficient policy opti-  
2 mization but is susceptible to distractions. We address this by developing Policy-  
3 Shaped Prediction (PSP), a method that empowers agents to interpret their own  
4 policies and shape their world models accordingly. By combining gradient-based  
5 interpretability, pretrained segmentation models, and adversarial learning, PSP  
6 outperforms existing distractor-reduction approaches. This work represents an  
7 interpretability-driven advance towards robust MBRL.

## 8 1 Introduction

9 Model-based reinforcement learning (MBRL) offers a promising path to data-efficient policy learning,  
10 demonstrating impressive performance with high-dimensional sensory data [Hafner et al., 2023].  
11 However, MBRL world models are particularly susceptible to distracting stimuli, a challenge that  
12 persists despite numerous attempts to address it [Deng et al., 2022, Fu et al., 2021, Wang et al., 2022,  
13 Seo et al., 2022, Wu et al., 2023, Schrittwieser et al., 2020].

14 We introduce Policy-Shaped Prediction (PSP), a novel method that uses gradient-based interpretability  
15 to identify and focus on important parts of an image-based environment. PSP interprets its own policy  
16 to prioritize relevant information, synergizing task-informed gradient-based loss weighting with a  
17 pre-trained segmentation model [Kirillov et al., 2023]. This approach creates a distraction-suppressing  
18 agent that outperforms leading image-based MBRL agents, particularly excelling against challenging  
19 and intricate, yet learnable, distractors. Our key contributions include:

- 20 • The development of PSP, combining gradient-based interpretability with pretrained segmen-  
21 tation to focus learning on important environment features.
- 22 • A challenging new benchmark for testing robustness to learnable distractions.
- 23 • Demonstration of PSP’s 2x improvement in robustness against challenging distractions  
24 while maintaining good performance in non-distracting settings.

## 25 2 Method

26 We introduce PSP, a method to reduce an agent’s sensitivity to useless distractions by focusing on  
27 sensory stimuli that are most relevant to its policy, rather than seeking to model everything in the  
28 environment. Our guiding intuition is that gradient-based interpretability techniques, traditionally  
29 used for post-hoc analysis, can be leveraged during training to highlight pixels in the environment  
30 that are important to the agent’s policy. Additionally, using image segmentation we aggregate these  
31 pixelwise salience signals to identify important objects.

32 PSP employs (1) gradients of the policy with respect to image inputs to identify task-relevant elements  
 33 of the image, (2) a segmentation model to aggregate gradients within each object in the image, and  
 34 (3) an adversarial objective to the image encoder of the world model that discourages encoding  
 35 of duplicate information about the previous action. Figure 1 illustrates the training modifications  
 36 made by this method to the underlying DreamerV3 [Hafner et al., 2023] architecture. Notably, since  
 37 these modifications only affect the training stage of the world model, the DreamerV3 agent remains  
 38 unaltered during inference. Below, we describe each of the three key components in detail.

### 39 2.1 Task-informed image reconstruction through interpretability-based weighting

40 Our approach builds upon the core idea that signals most important to the actor and/or critic should  
 41 be given special importance in the world model. We extend the concepts of Value-Gradient weighted  
 42 Model loss (VaGraM) [Voelcker et al., 2022] to high dimensional image inputs, which the previous  
 43 work did not demonstrate. This extension to the image domain is inspired by gradient-based  
 44 interpretability methods such as saliency maps [Simonyan et al., 2013, Shrikumar et al., 2017,  
 45 Ancona et al., 2019]. Doing so requires novel work mitigating the problems of using gradient signals  
 46 for high dimensional image input rather than low dimensional proprioceptive input. Additionally, we  
 47 test with complex signals that are present in the same image inputs that contain useful information,  
 48 whereas VaGraM tests on simple additional appended "distractor dimensions", which are independent  
 49 of the state space and reward function.

50 While VaGraM focused solely on the using a gradient signal from the value function, we hypothesize  
 51 that the gradient of the policy may provide an even more informative signal. To compute the policy-  
 52 gradient weighting, we first sum across the dimensions of the action vector  $\mathbf{a} = \mathbb{E}(\pi(\mathbf{s}))$ , where  $\mathbf{s}$   
 53 is the latent state of the world model, to produce a scalar  $a = \sum_j a_j$ , and then take the gradient  
 54 with respect to the pixels of the input image  $x$ . To apply this weighting in the context of DreamerV3  
 55 [Hafner et al., 2023], we scale the image reconstruction loss term at each pixel  $i$ , for reconstructed  
 56 image  $\hat{x}_i$ .

$$\mathcal{L}_{\text{image}}(\phi) = \sum_{x_i} \frac{\partial a}{\partial x_i} (\hat{x}_i - x_i)^2 \quad (1)$$

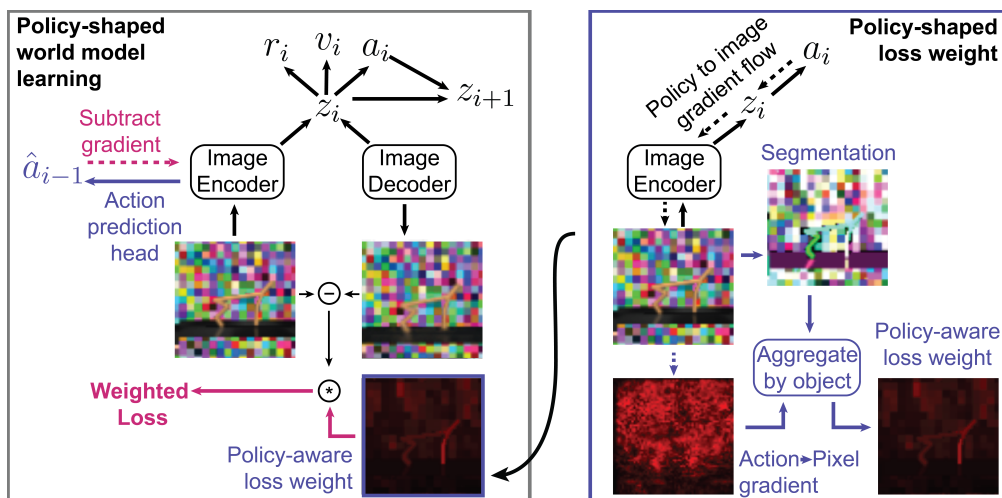


Figure 1: Policy-Shaped Prediction in an environment with challenging distractions. (left) Training of an otherwise-unaltered DreamerV3 agent is modified in two ways: 1) A head is added to predict the previous action based on the image encoding, and the gradient of the head is subtracted from the gradient of the image encoder, and 2) the loss is scaled pixelwise by a policy-shaped loss weight. (right) The loss weight uses the gradient of the policy to the input pixels. The image is segmented, and the pixel weights are averaged within each segmented object. Dashed lines signify gradient flow.

## 57 2.2 Improving saliency maps with object-based aggregation of gradient weights

58 Gradient-based weighting of the world model’s reconstruction faces challenges due to noisiness from  
59 small-scale fluctuations [Smilkov et al., 2017]. While more computationally demanding approaches  
60 exist [Sundararajan et al., 2017, Smilkov et al., 2017], we introduce a novel, efficient solution:  
61 object-based aggregation of explainability signals using the objects detected by any high-quality  
62 segmentation model (SEG). We used the Segment Anything Model (SAM) [Kirillov et al., 2023], but  
63 other models of sufficient quality may be utilized.

64 During data collection, we segment each image into object masks, including a mask for unassigned  
65 pixels. The weight of a pixel  $x_i$  in segment  $\text{SEG}(x_i)$  is:

$$W_i = \frac{1}{|\text{SEG}(x_i)|} \sum_{j \in \text{SEG}(x_i)} |\partial a / \partial x_j| \quad (2)$$

66 We clip the raw saliency map to the 99th percentile before aggregation. If all gradients are zero,  
67 we set  $W_i = 1$  for all  $i$ . We also linearly interpolate between the saliency weighting and a uniform  
68 weighting:  $W_i'' = \alpha W_i' + (1 - \alpha)$  where  $W_i' = \text{width} \cdot \text{height} \cdot W_i / \sum_i W_i$  and  $\alpha = 0.9$ . This allows  
69 the world model to maintain reasonable reconstruction of less-salient aspects of the environment.

## 70 2.3 Adversarial action prediction head

71 The DreamerV3 world model consists of three main components: a convolutional neural network  
72 (CNN) image encoder  $z_t \sim q_\phi(z_t|h_t, e_t)$  with  $e_t = \text{CNN}_\rho(x_t)$ , which processes the input image,  
73 serves as a prior during training, and encodes the environment state during inference; a recurrent  
74 state space machine (RSSM) consisting of  $h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1})$  and  $\hat{z}_t \sim p_\phi(\hat{z}_t|h_t)$  that  
75 is trained to simulate the progression of latent states given actions; and an image decoder,  $\hat{x}_t \sim$   
76  $p_\phi(\hat{x}_t|h_t, z_t)$  which reconstructs the image from the latent state. Problematically, the encoder can  
77 capture information about previous actions from the image, despite this information already being  
78 provided directly to the RSSM through the action input. In other words,  $z_t$  may source information  
79 about  $a_{t-1}$  directly through  $x_t$ , despite  $a_{t-1}$  being an argument to  $f_\phi$  during the computation of  
80  $h_t$ . Unfortunately, our reconstruction loss weighting may not solve this problem, since during  
81 backpropagation from the actor-critic functions, we do not distinguish information about previous  
82 actions that comes from the image versus the action input to the RSSM.

83 To prevent the encoder from redundantly capturing previous action information already provided to  
84 the RSSM, we introduce an adversarial MLP head:

$$\hat{a}_{t-1} = \text{MLP}_\omega(\text{CNN}_\rho(\text{sg}(x_t))) \quad (3)$$

$$\mathcal{L}_{\text{AdvHead}}(\hat{a}_{t-1}, a_{t-1}) = (\hat{a}_{t-1} - a_{t-1})^2 \quad (4)$$

85 During world model training, we subtract the scaled gradient  $\epsilon \cdot \nabla_\theta \mathcal{L}(\hat{a}_{t-1}, a_{t-1})$  from the overall  
86 gradient ( $\epsilon = 1e3$ ), ensuring action information comes solely from the provided action vector.

87 Our training procedure for a DreamerV3 agent is shown in Algorithm 1. We note that it should be  
88 possible to apply these concepts of gradient-based weighting, segmentation-based aggregation, and  
89 adversarial action prediction to world models other than our chosen DreamerV3 architecture.

## 90 3 Experiments

91 To evaluate the model’s performance we design our experiments around the following questions:

- 92 Q1. Is our agent robust against distractors which are learnable by the world model, but of no  
93 utility for the actor-critic?
- 94 Q2. What aspects of the environment are assigned importance by our method?
- 95 Q3. Is our agent robust against distractors that are unrelated to the agent’s actions?
- 96 Q4. Does our agent maintain performance in standard, lower-distraction environments?
- 97 Q5. What are the contributions of each component of our method?

---

**Algorithm 1** Policy-Shaped Prediction training (for DreamerV3)

---

```
1: Input: World model parameterized by  $\phi$ , policy  $\pi$  parameterized by  $\theta$ , image encoder
   parameterized by  $\rho$ , replay buffer with image transitions  $(x_{t-1}, a_{t-1}, x_t, r_t, c_t)$ , SEG segmentation
   model (SAM, in our application), action prediction MLP parameterized by  $\omega$ 
2: for training iteration 1, 2, ... do
3:   Sample batch of transition sequences
4:    $G = \nabla_x \pi_\theta$  # Gradient of policy with respect to input image pixels
5:    $S = \text{SEG}(x)$  # Segmentation of input image
6:    $W = \text{agg}(G, S)$  # Aggregate gradient using segmentation
7:    $W'_i = W_i / \sum_i W_i$  # Normalize weighting
8:    $W'' = \alpha W' + (1 - \alpha) \mathbf{1}_{\text{shape}}(W')$  # Linearly interpolate with uniform weighting
9:    $\mathcal{L}_{\text{pred}}(\phi) = -\ln p_\phi(x_t | z_t, h_t) \odot W'' - \ln p_\phi(r_t | z_t, h_t) - \ln p_\phi(c_t | z_t, h_t)$ 
   # Weighted DreamerV3 prediction loss
10:   $\mathcal{L}(\phi) = \mathbb{E}_{q_\phi} \left[ \sum_{t=1}^T (\beta_{\text{pred}} \mathcal{L}_{\text{pred}}(\phi) + \beta_{\text{dyn}} \mathcal{L}_{\text{dyn}}(\phi) + \beta_{\text{rep}} \mathcal{L}_{\text{rep}}(\phi)) \right]$  # DreamerV3 model loss
11:   $\hat{a}_{t-1} = \text{MLP}_\omega(\text{stop\_gradient}(\text{CNN}_\rho(x_t)))$  # Adversarial action prediction head
12:   $\phi \leftarrow \text{Adam}(\nabla \mathcal{L} - \epsilon * \partial \mathcal{L}(\hat{a}_{t-1}, a_{t-1}) / \partial \rho, \phi)$ 
13:   $\mathcal{L}_{\text{AdvHead}}(\hat{a}_{t-1}, a_{t-1}) = (\hat{a}_{t-1} - a_{t-1})^2$ 
14:   $\omega \leftarrow \text{Adam}(\nabla \mathcal{L}_{\text{AdvHead}}, \omega)$ 
15: end for
```

---

### 98 3.1 Experimental details

99 **Baselines** We test four Model-Based RL approaches as baselines: DreamerV3 [Hafner et al., 2023],  
100 and three methods specifically designed to handle distractions – Task Informed Abstractions [Fu  
101 et al., 2021], Denoised MDP (method in Figure 2b) [Wang et al., 2022], and DreamerPro [Deng et al.,  
102 2022]. Additionally, we choose DrQv2 [Yarats et al., 2021a] as a representative baseline Model-Free  
103 approach. For all agents, we use 3 random seeds per task, and default hyperparameters.

104 **Environment details** Visual observations are  $64 \times 64 \times 3$  pixel renderings. We test performance in  
105 three environments: DeepMind Control Suite (DMC) [Tassa et al., 2018], Reafferent DMC (described  
106 below), and Distracting Control Suite [Stone et al., 2021] (background video initialized to a random  
107 frame each episode, 2,000 grayscale frames from the "driving car" Kinetics dataset [Kay et al., 2017]).  
108 For each environment, we test two tasks: Cheetah Run and Hopper Stand. We selected these tasks  
109 because they present different levels of difficulty, allowing us to assess how distraction-sensitivity  
110 depends on the inherent difficulty of the task. For ablation experiments, we test on Cheetah Run.

### 111 3.2 Reafferent Deepmind Control Suite

112 To test our hypothesis, we devised the Reafferent  
113 Deepmind Control environment, inspired by  
114 [Stone et al., 2021]. This environment features  
115 distracting backgrounds that depend deterministically  
116 on the agent’s previous action and elapsed  
117 time, mimicking complex self-generated distractors  
118 in the natural world. The background consists of  
119 2,500  $16 \times 16$  color grids, mapped to 625  
120 time values and 4 discretized values of the first  
121 action dimension.

122 Many methods encode assumptions about the  
123 forms distractors will take (usually uncorrelated  
124 to agent actions, reward, or both), rather than a  
125 means of generally identifying and ignoring dis-  
126 tractors. We hypothesize that a learning-based  
127 approach, in which we avoid distraction by learn-  
128 ing what is actually important for the agent to  
129 get things done, has the potential to overcome  
130 even learnable-but-not-useful distractions.

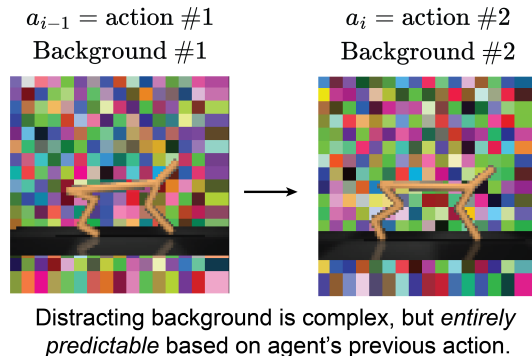


Figure 2: Schematic of the Reafferent Deepmind Control environment. The distracting background is entirely predictable based on the agent’s previous action and the elapsed time in the episode.

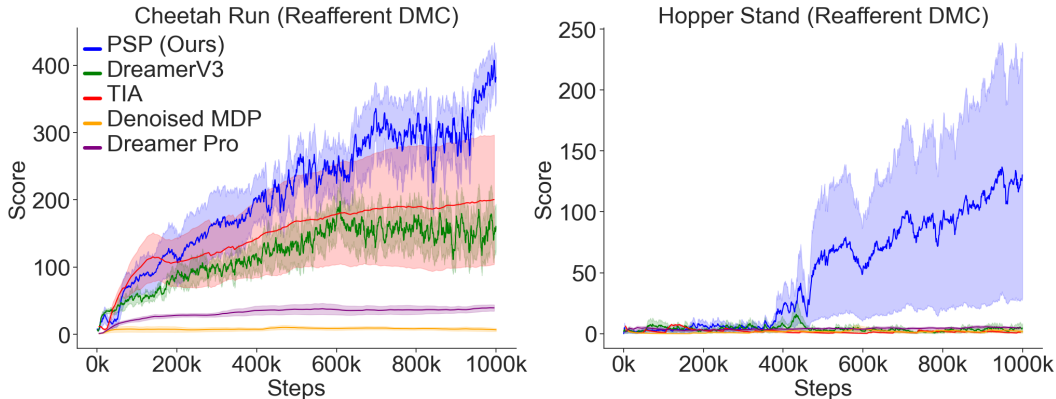


Figure 3: Training curve comparisons on Reafferent Deepmind Control. Mean  $\pm$  std. err.

131 We find that baseline MBRL methods perform poorly in this environment (Table 1, Figure 3), often reproducing the distracting background at the expense of accurately modeling the agent (Figure 132 4). In contrast, our method demonstrates substantial improvement 133 over existing baselines, achieving scores beyond their reach despite 134 some variance in performance and **affirmatively answering Q1**. 135 136

137 The model-free DrQv2 agent demonstrates robust performance, as 138 expected since its CNN encoder is learned as part of the policy. 139 In contrast, model-based methods face challenges when the world 140 model’s learning objective differs from the policy’s. Our method 141 bridges this gap, achieving superior performance while retaining the 142 advantages of model-based RL.

143 PSP demonstrates a substantial improvement 144 over the baselines (Table 1, Figure 3). Although 145 it shows a higher than desired level of variance 146 between runs, especially on the more challenging 147 Hopper Stand task, it nevertheless achieves 148 scores beyond the reach of any of the baselines. 149 We note that none of the 12 runs across the 4 150 baseline methods demonstrate a score substantially 151 above 0.

### 152 3.3 Performance on unaltered 153 DMC and Distracting Control Suite

154 Importantly, PSP performs comparably to other 155 methods (including DreamerV3) on the unaltered 156 Deepmind Control Suite, demonstrating

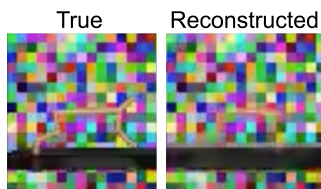


Figure 4: Denoised MDP reconstructs the background with a high degree of fidelity, but does not clearly render the Cheetah agent.

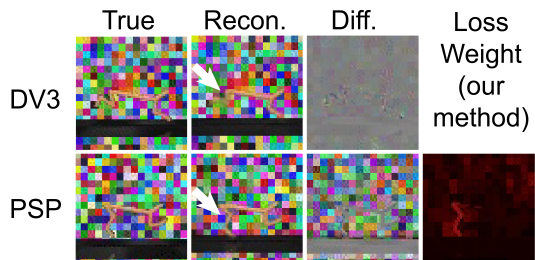


Figure 5: PSP vs. DreamerV3 on Reafferent Cheetah Run. From left: true, reconstructed, difference (true - recon.), loss weighting of PSP. DreamerV3 reproduces the background but not the back leg (see white arrow), and PCP renders the leg while not bothering to accurately model the background, **answering the question posed in Q2**.

Table 1: Performance comparison across environments. DrQv2 is model-free, all others are model-based. TIA is task-informed abstraction, dMDP is denoised MDP, mean  $\pm$  standard deviation.

Task	DrQv2	DreamerV3	DreamerPro	TIA	dMDP	PSP
Reafferent Control						
Cheetah Run	565.1 $\pm$ 35.5	158.4 $\pm$ 45.7	39.7 $\pm$ 9.0	200.4 $\pm$ 203.9	6.7 $\pm$ 4.3	<b>383.1 <math>\pm</math> 23.8</b>
Hopper Stand	210.3 $\pm$ 353.8	4.6 $\pm$ 3.9	3.8 $\pm$ 1.0	0.9 $\pm$ 0.3	1.7 $\pm$ 2.5	<b>128.5 <math>\pm</math> 215.7</b>
Unmodified Deepmind Control						
Cheetah Run	736.0 $\pm$ 17.0	521.1 $\pm$ 136.3	908.4 $\pm$ 1.6	773.7 $\pm$ 22.7	763.0 $\pm$ 62.8	712.3 $\pm$ 32.3
Hopper Stand	752.9 $\pm$ 206.8	867.4 $\pm$ 15.9	890 $\pm$ 11.2	298.4 $\pm$ 512	897.9 $\pm$ 14.2	865.6 $\pm$ 53.6
Distracting Deepmind Control						
Cheetah Run	364.4 $\pm$ 60.7	243.8 $\pm$ 81.2	179.1 $\pm$ 24	548.5 $\pm$ 238.9	397.4 $\pm$ 111.8	408.6 $\pm$ 125.1
Hopper Stand	781.1 $\pm$ 110.3	173.7 $\pm$ 160.9	561.8 $\pm$ 103.1	200.5 $\pm$ 171.7	13.2 $\pm$ 16.5	417.7 $\pm$ 118.9

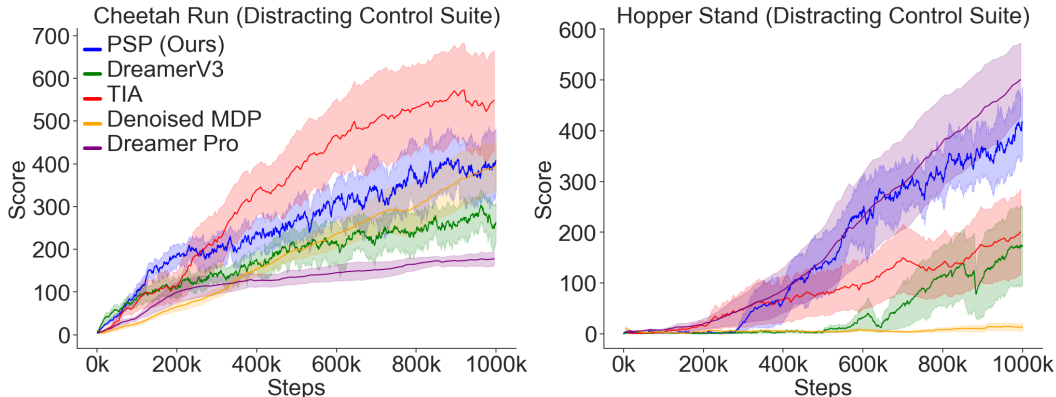


Figure 6: Training curve comparison on Distracting Control. Mean  $\pm$  std. err.

Table 2: Performance of ablated versions of PSP (for refferent and unaltered Cheetah Run).

Gradient weighting	Gradient weighting with segmentation	Unaltered	Refferent
Policy	✓	712.3 $\pm$ 32.3	383.1 $\pm$ 23.8
Policy	✗	742.1 $\pm$ 79.7	188.4 $\pm$ 9.4
None	✗	521.1 $\pm$ 136.3	158.4 $\pm$ 45.7

157 that we have not introduced a tradeoff between  
 158 performance on distracting and non-distracting environments and **addressing Q4**. (Table 1, Figure  
 159 A1).

160 On Distracting Control tasks, in which the background distractor is uncoupled from the agent’s  
 161 actions, PSP produced consistently improved performance relative to baseline DreamerV3, in contrast to  
 162 the more variable performance of DreamerPro, TIA, and Denoised MDP (Table 1, Figure 6),  
 163 **addressing Q3**.

164 In sum, PSP exhibits similar performance to  
 165 baseline methods in commonly used tests of  
 166 distractor-suppression and in non-distracting  
 167 environments, while also demonstrating un-  
 168 matched performance on particularly challeng-  
 169 ing distractors that are complex but learnable.  
 170 Given the success of MBRL in non-adversarial  
 171 environments, even when compared with lead-  
 172 ing Model Free Reinforcement Learning tech-  
 173 niques [Hafner et al., 2023], this work points to  
 174 ways of matching these gains in an adversarial  
 175 setting.

### 176 3.4 Ablation study

177 To understand the contributions of each sub-  
 178 component of the method (**Q5**), we conduct ab-  
 179 lations on the refferent and unaltered Cheetah Run (Table 2). We find that while some ablations  
 180 trade off performance between the environments, our complete model has good performance on  
 181 both. In particular, segmentation-based aggregation is critical to improving our model’s performance  
 182 amid distractors, while also maintaining its performance in the non-adversarial baseline. Overall,  
 183 the results of the ablations confirm that combining segmentation, policy gradient sensory weight-  
 184 ing, and adversarial action prediction results in the best scores across the unaltered and refferent  
 185 environments.

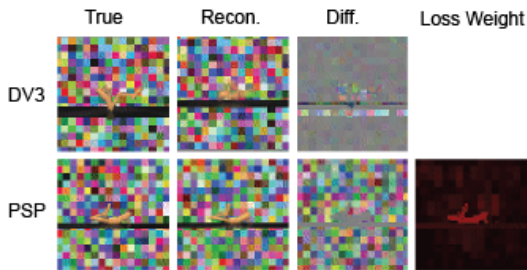


Figure 7: PSP vs. DreamerV3 on Refferent Hopper Stand. From left: true, reconstructed, difference (true - recon.), loss weighting of PSP. DreamerV3 reproduces the background but not the agent, and PSP renders the agent while not bothering to accurately model the background, **answering the question posed in Q2**.

## 186 4 Related Work

187 **Related Work** Recent advances in Model Based RL (MBRL) like World Models [Ha and Schmid-  
188 huber, 2018], SimPLe [Kaiser et al., 2019], MuZero [Schrittwieser et al., 2020], EfficientZero [Ye  
189 et al., 2021], and DreamerV3 [Hafner et al., 2023] have shown impressive performance but remain sus-  
190 ceptible to distractions [Lambert et al., 2020]. Various approaches have been proposed to address this,  
191 including (1) structural regularizations (DreamerPro [Deng et al., 2022], Agent Control-Endogenous  
192 State Discovery [Lamb et al., 2022], Task Informed Abstractions [Fu et al., 2021], Denoised MDPs  
193 [Wang et al., 2022]), ensemble methods [Clavera et al., 2018], and (2) learning-based approaches  
194 that use actor-critic functions to guide world modeling (VaGraM [Voelcker et al., 2022], Mismatched  
195 No More Eysenbach et al. [2022], Goal-Aware Prediction [Nair et al., 2020], Masked world models  
196 for visual control [Seo et al., 2023], The value equivalence principle for model-based reinforcement  
197 learning [Grimm et al., 2020], MuZero [Schrittwieser et al., 2020], Value Prediction Networks [Oh  
198 et al., 2017]). Parallel work in Model Free RL (MFRL) has also tackled distraction sensitivity, with  
199 methods like DrQv2 [Yarats et al., 2021a] and approaches using attention mechanisms [Mott et al.,  
200 2019], prototypes [Yarats et al., 2021b], and dynamic sparse training [Grooten et al., 2023a,b].

## 201 5 Discussion

202 PSP introduces a novel approach to model-based reinforcement learning that leverages interpretability  
203 techniques not just for analysis, but as an integral part of the learning process. By allowing the agent  
204 to interpret its own policy, PSP focuses the world model’s capacity on aspects of the environment  
205 most relevant for decision making. This self-interpretation process comprises three key components:  
206 1) We use gradient-based interpretability methods, analogous to saliency maps [Simonyan et al.,  
207 2013], to identify important pixels in the input image; 2) We aggregate pixel importance by object  
208 using a pre-trained segmentation model, providing a higher-level interpretation of the environment;  
209 3) We employ an adversarial prediction head to prevent wasteful encoding of known information.

210 Our work opens avenues for future research in interpretable RL, such as using more advanced  
211 explainability gradient-based attribution methods like Integrated Gradients [Sundararajan et al., 2017,  
212 Ancona et al., 2019]. While PSP demonstrates promising results, it has limitations, including its  
213 object-centric assumptions and the computational requirements of the segmentation model.

214 **Outlook** In conclusion, PSP represents a significant step towards robust model-based RL via the  
215 direct integration of of model interpretability techniques. The findings here open other lines of inquiry  
216 such as using more explainable architectures, utilizing faster segmentation models and utilizing  
217 segmentation models designed for videos in order to do temporal aggregation.

## 218 References

- 219 M. Ancona, E. Ceolini, C. Öztireli, and M. Gross. Gradient-based attribution methods. *Explainable*  
220 *AI: Interpreting, explaining and visualizing deep learning*, pages 169–191, 2019.
- 221 I. Clavera, J. Rothfuss, J. Schulman, Y. Fujita, T. Asfour, and P. Abbeel. Model-based reinforcement  
222 learning via meta-policy optimization, 2018.
- 223 F. Deng, I. Jang, and S. Ahn. Dreamerpro: Reconstruction-free model-based reinforcement learning  
224 with prototypical representations. In *International Conference on Machine Learning*, pages  
225 4956–4975. PMLR, 2022.
- 226 B. Eysenbach, A. Khazatsky, S. Levine, and R. R. Salakhutdinov. Mismatched no more: Joint model-  
227 policy optimization for model-based rl. *Advances in Neural Information Processing Systems*, 35:  
228 23230–23243, 2022.
- 229 X. Fu, G. Yang, P. Agrawal, and T. Jaakkola. Learning task informed abstractions. In *International*  
230 *Conference on Machine Learning*, pages 3480–3491. PMLR, 2021.
- 231 C. Grimm, A. Barreto, S. Singh, and D. Silver. The value equivalence principle for model-based  
232 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:5541–5552, 2020.

- 233 B. Grooten, G. Sokar, S. Dohare, E. Mocanu, M. E. Taylor, M. Pechenizkiy, and D. C. Mocanu.  
 234 Automatic noise filtering with dynamic sparse training in deep reinforcement learning. *arXiv*  
 235 *preprint arXiv:2302.06548*, 2023a.
- 236 B. Grooten, T. Tomilin, G. Vasan, M. E. Taylor, A. R. Mahmood, M. Fang, M. Pechenizkiy, and D. C.  
 237 Mocanu. Madi: Learning to mask distractions for generalization in visual deep reinforcement  
 238 learning. *arXiv preprint arXiv:2312.15339*, 2023b.
- 239 D. Ha and J. Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 240 D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. Mastering diverse domains through world models.  
 241 *arXiv preprint arXiv:2301.04104*, 2023.
- 242 L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn,  
 243 P. Kozakowski, S. Levine, et al. Model-based reinforcement learning for atari. *arXiv preprint*  
 244 *arXiv:1903.00374*, 2019.
- 245 W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green,  
 246 T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The kinetics human action video dataset.  
 247 *CoRR*, abs/1705.06950, 2017. URL <http://arxiv.org/abs/1705.06950>.
- 248 A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg,  
 249 W.-Y. Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on*  
 250 *Computer Vision*, pages 4015–4026, 2023.
- 251 A. Lamb, R. Islam, Y. Efroni, A. Didolkar, D. Misra, D. Foster, L. Molu, R. Chari, A. Krishnamurthy,  
 252 and J. Langford. Guaranteed discovery of control-endogenous latent states with multi-step inverse  
 253 models. *arXiv preprint arXiv:2207.08229*, 2022.
- 254 N. Lambert, B. Amos, O. Yadan, and R. Calandra. Objective mismatch in model-based reinforcement  
 255 learning. *arXiv preprint arXiv:2002.04523*, 2020.
- 256 A. Mott, D. Zoran, M. Chrzanowski, D. Wierstra, and D. Jimenez Rezende. Towards interpretable re-  
 257 inforcement learning using attention augmented agents. *Advances in neural information processing*  
 258 *systems*, 32, 2019.
- 259 S. Nair, S. Savarese, and C. Finn. Goal-aware prediction: Learning to model what matters. In  
 260 *International Conference on Machine Learning*, pages 7207–7219. PMLR, 2020.
- 261 J. Oh, S. Singh, and H. Lee. Value prediction network. *Advances in neural information processing*  
 262 *systems*, 30, 2017.
- 263 J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart,  
 264 D. Hassabis, T. Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned  
 265 model. *Nature*, 588(7839):604–609, 2020.
- 266 Y. Seo, K. Lee, S. James, and P. Abbeel. Reinforcement learning with action-free pre-training from  
 267 videos, 2022.
- 268 Y. Seo, D. Hafner, H. Liu, F. Liu, S. James, K. Lee, and P. Abbeel. Masked world models for visual  
 269 control. In *Conference on Robot Learning*, pages 1332–1344. PMLR, 2023.
- 270 A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating  
 271 activation differences. In *International conference on machine learning*, pages 3145–3153. PMLR,  
 272 2017.
- 273 K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image  
 274 classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- 275 D. Smilkov, N. Thorat, B. Kim, F. B. Viégas, and M. Wattenberg. Smoothgrad: removing noise by  
 276 adding noise. *CoRR*, abs/1706.03825, 2017. URL <http://arxiv.org/abs/1706.03825>.
- 277 A. Stone, O. Ramirez, K. Konolige, and R. Jonschkowski. The distracting control suite—a challenging  
 278 benchmark for reinforcement learning from pixels. *arXiv preprint arXiv:2101.02722*, 2021.



- 279 M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *International*  
280 *conference on machine learning*, pages 3319–3328. PMLR, 2017.
- 281 Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel,  
282 A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- 283 C. Voelcker, V. Liao, A. Garg, and A.-m. Farahmand. Value gradient weighted model-based rein-  
284 forcement learning. *arXiv preprint arXiv:2204.01464*, 2022.
- 285 T. Wang, S. S. Du, A. Torralba, P. Isola, A. Zhang, and Y. Tian. Denoised mdps: Learning world  
286 models better than the world itself. *arXiv preprint arXiv:2206.15477*, 2022.
- 287 J. Wu, H. Ma, C. Deng, and M. Long. Pre-training contextualized world models with in-the-wild  
288 videos for reinforcement learning, 2023.
- 289 D. Yarats, R. Fergus, A. Lazaric, and L. Pinto. Mastering visual continuous control: Improved  
290 data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021a.
- 291 D. Yarats, R. Fergus, A. Lazaric, and L. Pinto. Reinforcement learning with prototypical rep-  
292 resentations. In *International Conference on Machine Learning*, pages 11920–11931. PMLR,  
293 2021b.
- 294 W. Ye, S. Liu, T. Kurutach, P. Abbeel, and Y. Gao. Mastering atari games with limited data. *Advances*  
295 *in neural information processing systems*, 34:25476–25488, 2021.

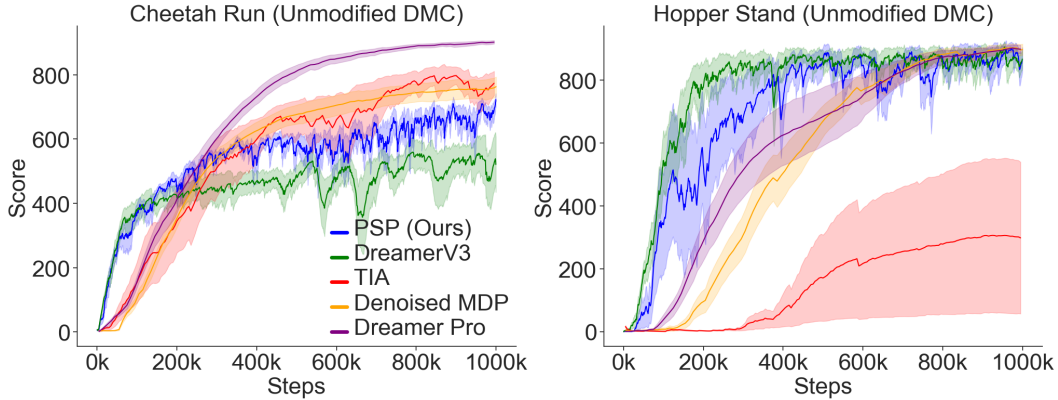


Figure A1: Training curve comparison on unmodified Deepmind Control. Mean  $\pm$  std. err.

## 296 A Broader Impacts

297 At the current stage, this work remains reasonably far from any large societal impacts, as it is limited  
 298 to agents interacting with small, simulated environments. Over the long term, however, if model-based  
 299 RL algorithms are used to control robots or internet-connected agents (such as large language model  
 300 agents), the potential for both large positive and negative societal impacts becomes relevant. On the  
 301 positive side, intelligent agents that are capable of modeling the world and avoiding distractors have  
 302 the potential to aid humans in a wide variety of scenarios, from housework, to medical applications,  
 303 to exploration, to internet research. On the negative side, agents without proper safeguards have the  
 304 potential to inflict harm on humans and the environment, whether through negligence or malfeasance.  
 305 Ultimately, our work is targeted at producing the positive impacts, while still allowing for mitigation  
 306 of the negative impacts.

## 307 B Experiments Compute Resources

308 Each trial of the PSP method used 4 Nvidia A40 GPUs to train the modified DreamerV3 model, and 4  
 309 A40 GPUs to run the segment anything model in parallel. Given an estimated 17 unique experiments  
 310 for the final paper, 3 trials per experiment with our method, and about 1.5 days per training run, we  
 311 used about  $17 * 3 * 1.5 * 8 \text{ GPUs} = 612 \text{ GPU days}$  on A40 accelerators. Early experiments with  
 312 this methodology likely used an additional 300. Baseline trials could be run on only a single A40  
 313 GPU or a desktop NVIDIA 2070 SUPER, usually in less than a day, and accounted for a comparably  
 314 negligible level of resources.

315 We believe this level of resource consumption could be easily reduced. The modifications to the  
 316 DreamerV3 model do not attempt to benchmark the most costly components. We suspect our method  
 317 of parallelizing the new backpropagation from the policy to the image could be optimized further  
 318 from its naive Jax implementation. Additionally, SAM could be supplanted by a more efficient  
 319 segmentation model. We focused on establishing the basic technique with SAM, but replacing it  
 320 should be the subject of future work.

## 321 C Code

322 An anonymized version of the code with instructions for reproducing these experiments will be  
 323 available for reviewers at this anonymous [GitHub Repository](#).

## 324 **NeurIPS Paper Checklist**

### 325 **1. Claims**

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

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

329 Justification: The abstract and introduction clearly state our contributions and claims, and  
330 match the details in the results section.

331 Guidelines:

- 332 • The answer NA means that the abstract and introduction do not include the claims  
333 made in the paper.
- 334 • The abstract and/or introduction should clearly state the claims made, including the  
335 contributions made in the paper and important assumptions and limitations. A No or  
336 NA answer to this question will not be perceived well by the reviewers.
- 337 • The claims made should match theoretical and experimental results, and reflect how  
338 much the results can be expected to generalize to other settings.
- 339 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
340 are not attained by the paper.

### 341 **2. Limitations**

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

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

344 Justification: Yes, the paper clearly points out what we believe to be core limitations and  
345 assumptions of our work, as well as present limitations that we do not believe are inherent  
346 to the method.

347 Guidelines:

- 348 • The answer NA means that the paper has no limitation while the answer No means that  
349 the paper has limitations, but those are not discussed in the paper.
- 350 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 351 • The paper should point out any strong assumptions and how robust the results are to  
352 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
353 model well-specification, asymptotic approximations only holding locally). The authors  
354 should reflect on how these assumptions might be violated in practice and what the  
355 implications would be.
- 356 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
357 only tested on a few datasets or with a few runs. In general, empirical results often  
358 depend on implicit assumptions, which should be articulated.
- 359 • The authors should reflect on the factors that influence the performance of the approach.  
360 For example, a facial recognition algorithm may perform poorly when image resolution  
361 is low or images are taken in low lighting. Or a speech-to-text system might not be  
362 used reliably to provide closed captions for online lectures because it fails to handle  
363 technical jargon.
- 364 • The authors should discuss the computational efficiency of the proposed algorithms  
365 and how they scale with dataset size.
- 366 • If applicable, the authors should discuss possible limitations of their approach to  
367 address problems of privacy and fairness.
- 368 • While the authors might fear that complete honesty about limitations might be used by  
369 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
370 limitations that aren't acknowledged in the paper. The authors should use their best  
371 judgment and recognize that individual actions in favor of transparency play an impor-  
372 tant role in developing norms that preserve the integrity of the community. Reviewers  
373 will be specifically instructed to not penalize honesty concerning limitations.

### 374 **3. Theory Assumptions and Proofs**

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

377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430

Answer: [NA]

Justification: The claims of this paper are tested empirically.

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: The paper itself should include enough detail to reproduce our results with the open source SAM and Dreamerv3 models. Because the implementation is not trivial, we will also release the GitHub repository in the camera ready version, which includes an implementation of the core algorithm for DreamerV3 and a shared implementation of the test environments for DreamerV3 and every baseline. We have not included it in the review version as the GitHub will identify the authors.

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

431 some way (e.g., to registered users), but it should be possible for other researchers  
432 to have some path to reproducing or verifying the results.

### 433 5. Open access to data and code

434 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
435 tions to faithfully reproduce the main experimental results, as described in supplemental  
436 material?

437 Answer: [Yes]

438 Justification: An anonymized version of the code will be available at the linked [GitHub](#)  
439 [Repository](#) for reviewers.

440 Guidelines:

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

### 460 6. Experimental Setting/Details

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

464 Answer: [Yes]

465 Justification: We include all relevant details of our benchmark and (where we diverge from  
466 the default) models.

467 Guidelines:

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

### 473 7. Experiment Statistical Significance

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

476 Answer: [Yes]

477 Justification: All experiments are performed with three trials and std. dev is reported.

478 Guidelines:

- 479 • The answer NA means that the paper does not include experiments.
- 480 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
481 dence intervals, or statistical significance tests, at least for the experiments that support  
482 the main claims of the paper.

- 483 • The factors of variability that the error bars are capturing should be clearly stated (for  
484 example, train/test split, initialization, random drawing of some parameter, or overall  
485 run with given experimental conditions).
- 486 • The method for calculating the error bars should be explained (closed form formula,  
487 call to a library function, bootstrap, etc.)
- 488 • The assumptions made should be given (e.g., Normally distributed errors).
- 489 • It should be clear whether the error bar is the standard deviation or the standard error  
490 of the mean.
- 491 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
492 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
493 of Normality of errors is not verified.
- 494 • For asymmetric distributions, the authors should be careful not to show in tables or  
495 figures symmetric error bars that would yield results that are out of range (e.g. negative  
496 error rates).
- 497 • If error bars are reported in tables or plots, The authors should explain in the text how  
498 they were calculated and reference the corresponding figures or tables in the text.

## 499 8. Experiments Compute Resources

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

503 Answer: [Yes]

504 Justification: We have provided a transparent and reasonable estimate of compute require-  
505 ments in the appendix.

506 Guidelines:

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

## 515 9. Code Of Ethics

516 Question: Does the research conducted in the paper conform, in every respect, with the  
517 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

518 Answer: [Yes]

519 Justification: The only new dataset involved in this work was generated by a small Python  
520 procedure and has no privacy risks or ethical concerns. The implemented model modifies an  
521 open source repo and has had no interaction with human subjects.

522 Guidelines:

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

## 528 10. Broader Impacts

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

531 Answer: [Yes]

532 Justification: This is foundational work for improving MBRL and any societal impacts are at  
533 least one order removed, but we have outlined the possible societal impacts of improved  
534 MBRL in general.

535 Guidelines:

- 536 • The answer NA means that there is no societal impact of the work performed.
- 537 • If the authors answer NA or No, they should explain why their work has no societal
- 538 impact or why the paper does not address societal impact.
- 539 • Examples of negative societal impacts include potential malicious or unintended uses
- 540 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 541 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 542 groups), privacy considerations, and security considerations.
- 543 • The conference expects that many papers will be foundational research and not tied
- 544 to particular applications, let alone deployments. However, if there is a direct path to
- 545 any negative applications, the authors should point it out. For example, it is legitimate
- 546 to point out that an improvement in the quality of generative models could be used to
- 547 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 548 that a generic algorithm for optimizing neural networks could enable people to train
- 549 models that generate Deepfakes faster.
- 550 • The authors should consider possible harms that could arise when the technology is
- 551 being used as intended and functioning correctly, harms that could arise when the
- 552 technology is being used as intended but gives incorrect results, and harms following
- 553 from (intentional or unintentional) misuse of the technology.
- 554 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 555 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 556 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 557 feedback over time, improving the efficiency and accessibility of ML).

558 **11. Safeguards**

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

562 Answer: [NA]

563 Justification: This paper describes a foundational change to MBRL and introduces no new  
564 datasets that pose a risk for misuse.

565 Guidelines:

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

576 **12. Licenses for existing assets**

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

580 Answer: [Yes]

581 Justification: All previous work is cited and no proprietary code has been used beyond what  
582 is allowed by its license.

583 Guidelines:

- 584 • The answer NA means that the paper does not use existing assets.
- 585 • The authors should cite the original paper that produced the code package or dataset.
- 586 • The authors should state which version of the asset is used and, if possible, include a  
587 URL.

- 588 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 589 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 590 service of that source should be provided.
- 591 • If assets are released, the license, copyright information, and terms of use in the
- 592 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)
- 593 has curated licenses for some datasets. Their licensing guide can help determine the
- 594 license of a dataset.
- 595 • For existing datasets that are re-packaged, both the original license and the license of
- 596 the derived asset (if it has changed) should be provided.
- 597 • If this information is not available online, the authors are encouraged to reach out to
- 598 the asset's creators.

### 599 13. New Assets

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

601 provided alongside the assets?

602 Answer: [Yes]

603 Justification: Source code will be provided before publication if paper is accepted.

604 Guidelines:

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

### 613 14. Crowdsourcing and Research with Human Subjects

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

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

616 well as details about compensation (if any)?

617 Answer: [NA]

618 Justification: No human subjects.

619 Guidelines:

- 620 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 621 human subjects.
- 622 • Including this information in the supplemental material is fine, but if the main contribu-
- 623 tion of the paper involves human subjects, then as much detail as possible should be
- 624 included in the main paper.
- 625 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 626 or other labor should be paid at least the minimum wage in the country of the data
- 627 collector.

### 628 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human

629 Subjects

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

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

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

633 institution) were obtained?

634 Answer: [NA]

635 Justification: No human subjects.

636 Guidelines:

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



639  
640  
641  
642  
643  
644  
645  
646

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