

# ON ROLLOUTS IN MODEL-BASED REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model-based reinforcement learning (MBRL) seeks to enhance data efficiency by learning a model of the environment and generating synthetic rollouts from it. However, accumulated model errors during these rollouts can distort the data distribution, negatively impacting policy learning and hindering long-term planning. Thus, the accumulation of model errors is a key bottleneck in current MBRL methods. We propose *Infoprop*, a model-based rollout mechanism that separates aleatoric from epistemic model uncertainty and reduces the influence of the latter on the data distribution. Further, *Infoprop* keeps track of accumulated model errors along a model rollout and provides termination criteria to limit data corruption. We demonstrate the capabilities of *Infoprop* in the *Infoprop-Dyna* algorithm, reporting state-of-the-art performance in Dyna-style MBRL on common MuJoCo benchmark tasks while substantially increasing rollout length and data quality.

## 1 INTRODUCTION

Reinforcement learning (RL) has emerged as a powerful framework for solving complex decision-making tasks like racing Vasco et al. (2024); Kaufmann et al. (2023) and gameplay OpenAI et al. (2019); Bi & D’Andrea (2024). However, when applying RL in real-world scenarios, a significant challenge is data inefficiency, which hinders the practicality of standard RL methods. Model-based reinforcement learning (MBRL) addresses this issue by learning an internal model of the environment Deisenroth & Rasmussen (2011); Chua et al. (2018); Janner et al. (2019); Hafner et al. (2020). By generating simulated interactions through model rollouts, MBRL can make informed decisions while substantially reducing the need for real-world data collection.

The quality of data from model-based rollouts is critical for MBRL performance. Model errors can distort the data distribution and hurt policy learning. Long-horizon planning is desirable, however, often infeasible as model errors accumulate over time. This effect is demonstrated in Figure 1. Even for a simple toy example (described in Appendix B), we see the data distribution of model-based rollouts under the state-of-the-art Trajectory Sampling (TS) Chua et al. (2018) scheme diverging quickly from the ground truth distribution of environment rollouts. Thus, data from TS rollouts can even be harmful to policy learning after a couple of time steps. This is largely because the TS mechanism does not explicitly address the effect of model errors on the propagated data distribution.

To tackle this challenge, we propose *Infoprop*, a novel model-based rollout mechanism that mitigates data distortion by addressing two key questions: *How to propagate?* and *When to stop?* We build our mechanism on explicitly leveraging the ability of common MBRL models to distinguish between aleatoric uncertainty due to process noise and epistemic uncertainty due to lack of data Lakshminarayanan et al. (2017); Becker & Neumann (2022). Making use of this property leads to substantially improved data consistency as depicted in Figure 1. In particular, we

- estimate and remove the stochasticity due to model error from the predictive distribution;
- formulate stopping criteria based on information loss to limit error accumulation; and
- demonstrate the potential of *Infoprop* as a direct plugin to standard MBRL methods using the example of Dyna-style MBRL. The resulting *Infoprop-Dyna* algorithm yields state-of-the-art performance in MBRL on common MuJoCo tasks, while substantially improving the data consistency of model-based rollouts and thus allowing for longer rollout horizons.

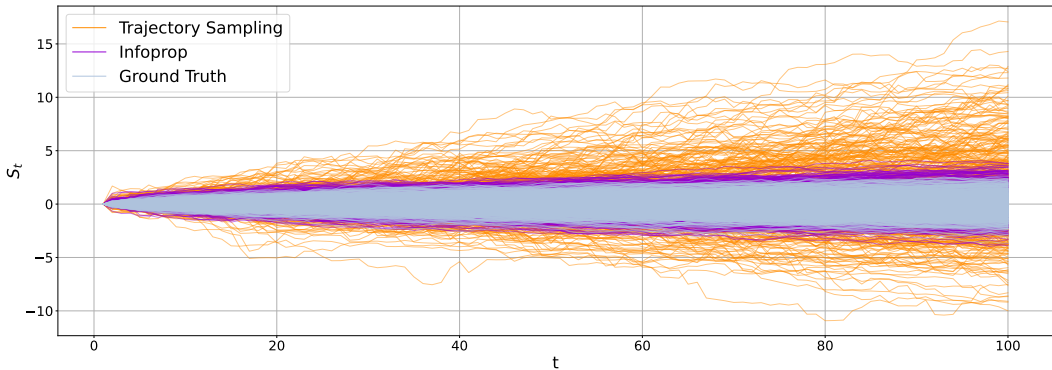


Figure 1: Comparing Data Consistency of Model-based Rollouts. *Trajectories under the proposed Infoprop mechanism follow the ground truth distribution of environment rollouts closely while rolling out the same model under the common TS scheme Chua et al. (2018) results in distorted data.*

## 2 BACKGROUND

In the following, we introduce the fundamental concepts of information theory and MBRL. Appendix A provides an overview of the notation introduced and used in the remainder of the paper.

### 2.1 INFORMATION THEORY

We will estimate the degree of data corruption in Infoprop rollouts using information-theoretic arguments. Information theory serves to quantify the uncertainty of a random variable (RV) Shannon (1948). Given the discrete RVs  $X : \Omega \rightarrow \mathcal{X}$  and  $Y : \Omega \rightarrow \mathcal{Y}$ , the marginal entropy  $\mathbb{H}(X) = -\sum_{x \in \mathcal{X}} \mathbb{P}[X = x] \log_2(\mathbb{P}[X = x])$  describes the average uncertainty about  $X$  in bits. Further, the conditional entropy  $\mathbb{H}(X|Y = y) = -\sum_{x \in \mathcal{X}} \mathbb{P}[X = x|Y = y] \log_2(\mathbb{P}[X = x|Y = y])$  gives the uncertainty about  $X$ , given a realization of  $Y$ . Based on marginal and conditional entropy, the reduction in uncertainty about  $X$  given a realization of  $Y$  is described by mutual information

$$\mathbb{I}(X; Y = y) = \mathbb{H}(X) - \mathbb{H}(X|Y = y), \quad (1)$$

with  $\mathbb{I}(X; Y = y) = 0$  if the RVs are independent. In the following, we focus on Gaussian RVs and use the notion of quantized entropy Cover & Thomas (2006) with details provided in Appendix D.1.

### 2.2 REINFORCEMENT LEARNING

Reinforcement learning addresses sequential decision-making problems where the environment is typically modeled as a discrete-time Markov decision process (MDP) represented by the tuple  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{R}, P_{\mathcal{R}}, P_{\mathcal{S}}, \xi_0, \gamma\}$ . Here,  $\mathcal{S} \subseteq \mathbb{R}^{n_S}$  denotes the state space with  $S_t \in \mathcal{S}$  being the RV of the state at time  $t$  and  $s_t$  its realization. Similarly,  $\mathcal{A} \subseteq \mathbb{R}^{n_A}$  represents the action space with  $A_t \in \mathcal{A}$  the RV and  $a_t$  the realization of the action as well as  $\mathcal{R} \subseteq \mathbb{R}$  the set of rewards with  $R_t \in \mathcal{R}$  and  $r_t$  the reward at time  $t$ . We make the common simplifying assumption Bellemare et al. (2023) that the next state and reward are independent given the current state-action pair. Thus, a transition step in the environment can be expressed concerning a reward kernel  $P_{\mathcal{R}}$  and a dynamics kernel  $P_{\mathcal{S}}$  as

$$R_{t+1} \sim P_{\mathcal{R}}(\cdot|S_t, A_t) \quad \text{and} \quad S_{t+1} \sim P_{\mathcal{S}}(\cdot|S_t, A_t). \quad (2)$$

Further, initial states are distributed according to  $S_0 \sim \xi_0$ , and actions according to the policy  $A_t \sim \pi(\cdot|S_t)$ . We aim to learn an optimal policy  $\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t R_{t+1}]$  that maximizes the expected sum of rewards discounted by  $\gamma \in [0, 1)$ , referred to as return.

### 2.3 MODEL-BASED REINFORCEMENT LEARNING

There are four main categories of MBRL that all build on model-based rollouts. (i) Dyna-style methods Sutton (1991); Janner et al. (2019) use model-based rollouts to generate training data for a model-free agent. (ii) Model-based planning approaches Chua et al. (2018); Williams et al. (2017);

108 Nagabandi et al. (2018); Hafner et al. (2019) do not learn an explicit policy but perform planning  
 109 via model rollouts during deployment. (iii) Analytic gradient methods Deisenroth & Rasmussen  
 110 (2011); Hafner et al. (2020; 2021; 2023) optimize the policy by backpropagating the performance  
 111 gradient through model rollouts. (iv) Value-expansion approaches Feinberg et al. (2018); Buckman  
 112 et al. (2018) stabilize the temporal difference target using model-based rollouts.

113 The model architecture of an MBRL algorithm determines the set of mechanisms for model rollouts.  
 114 In this work, we focus on rolling out the particularly successful class of aleatoric epistemic separator  
 115 (AES) models Lakshminarayanan et al. (2017); Becker & Neumann (2022) that can distinguish  
 116 aleatoric uncertainty corresponding to the estimate of process noise from epistemic uncertainty.  
 117

## 118 2.4 ENVIRONMENT INTERACTION VS. MODEL-BASED ROLLOUTS

120 Model-based rollouts aim to substitute environment interaction in MBRL. Thus, we compare the  
 121 data generation process through environment interaction to the process of model-based rollouts.

122 We model environment dynamics as a nonlinear function  $\mu(S_t, A_t)$  with additive heteroscedastic  
 123 process noise that is normally distributed with variance  $\Sigma(S_t, A_t)$ . Thus, environment rollouts, as  
 124 depicted in Figure 1, are generated by iterating the dynamics  
 125

$$126 S_{t+1} = \mu(S_t, A_t) + L(S_t, A_t)W_t, \quad (3)$$

127 with  $L(S_t, A_t)L(S_t, A_t)^\top = \Sigma(S_t, A_t)$  and the process noise  $W_t \sim \mathcal{N}(0, I)$ . Consequently, the  
 128 transition kernel<sup>1</sup> of the environment is defined as  $P_S(\cdot|S_t, A_t) = \mathcal{N}(\mu(S_t, A_t), \Sigma(S_t, A_t))$ .

130 In MBRL, however, we do not have access to  $P_S$  directly but typically rely on a parametric model  
 131 with the random parameters  $\Theta_t \in \vartheta$ . Besides estimates of nonlinear dynamics  $\hat{\mu}_{\Theta_t}(S_t, A_t)$  and  
 132 process noise  $\hat{\Sigma}_{\Theta_t}(S_t, A_t)$ , AES models provide an estimate of the parameter distribution  $\Theta_t \sim \mathbb{P}_\Theta$ ,  
 133 e.g. via ensembling Lakshminarayanan et al. (2017) or dropout Becker & Neumann (2022). These  
 134 models are typically propagated using the TS Chua et al. (2018) rollout mechanism via iterating  
 135

$$136 S_{t+1} = \hat{\mu}_{\Theta_t}(S_t, A_t) + \hat{L}_{\Theta_t}(S_t, A_t)W_t \quad (4)$$

137 with  $\hat{L}_{\Theta_t}(S_t, A_t)\hat{L}_{\Theta_t}(S_t, A_t)^\top = \hat{\Sigma}_{\Theta_t}(S_t, A_t)$ ,  $W_t \sim \mathcal{N}(0, I)$ , and  $\Theta_t \sim \mathbb{P}_\Theta$ . This results in the  
 138 TS rollouts in Figure 1 and induces the kernel  $\hat{P}_{S, \text{TS}}(\cdot|S_t, A_t) = \mathcal{N}(\hat{\mu}_{\Theta_t}(S_t, A_t), \hat{\Sigma}_{\Theta_t}(S_t, A_t))$ .  
 139 The majority of recent MBRL approaches use the TS rollout mechanism, e.g. Chua et al. (2018);  
 140 Becker & Neumann (2022); Janner et al. (2019); Pan et al. (2020); Yu et al. (2020); Luis et al.  
 141 (2023). Pseudocode is provided in Algorithm 2 of Appendix C.  
 142  
 143  
 144

## 145 3 PROBLEM STATEMENT

147 Revisiting Figure 1 allows us to illustrate the effects of different sources of stochasticity by compar-  
 148 ing environment interaction under  $P_S$  to TS rollouts under  $\hat{P}_{S, \text{TS}}$ . While different realizations of  
 149 process noise  $w_t \sim \mathcal{N}(0, I)$  allow for keeping track of the environment distribution, the sampling  
 150 process  $\theta_t \sim \mathbb{P}_\Theta$  introduces additional stochasticity that leads to an overestimated total variance in  
 151 the TS rollout distribution. This effect is amplified through the continued propagation of erroneous  
 152 predictions making data at later steps unfit for policy learning. We ask the following questions:  
 153

- 154 (i) How can we construct a predictive distribution closely resembling environment dynamics?
- 155 (ii) How can we quantify the degree of data corruption due to model error?
- 156 (iii) When should model-based rollouts be terminated due to data corruption?

157 We address these questions by proposing the *Infoprop* rollout mechanism. Infoprop isolates and  
 158 removes epistemic uncertainty for an improved predictive distribution, keeps track of data corruption  
 159 using information-theoretic arguments, and terminates rollouts based on the degree of corruption.  
 160

161 <sup>1</sup>As  $P_{\mathcal{R}}$  typically is a known deterministic function in the context of MBRL, while  $P_S$  is the unknown  
 object we aim to model, the discussion henceforth focuses on approximating  $P_S$  without loss of generality.

## 4 INFOPROP ROLLOUT MECHANISM

In the following, we introduce the Infoprop mechanism for model-based rollouts. As depicted in Figure 2, we decompose model predictions into a signal fraction representing the environment dynamics and noise fraction introduced by model error. This perspective allows to interpret model rollouts as communication through a noisy channel. We estimate both the signal and the noise distribution and use these to infer a belief over the environment state, given an observation of the model state. This belief state represents the foundation of the Infoprop rollout mechanism.

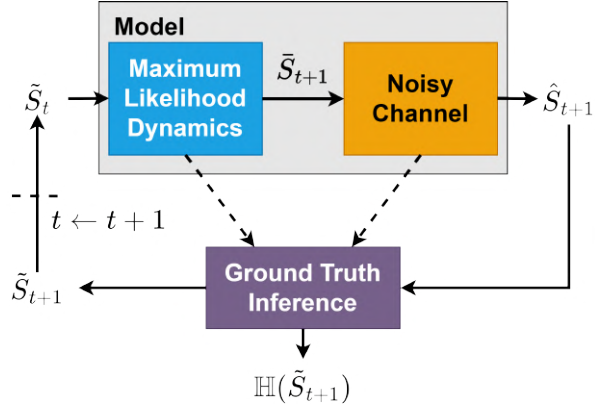


Figure 2: Infoprop block diagram

### 4.1 THEORETICAL SETUP

First, we introduce additional notation to specify RVs under different transition kernels.

**Definition 1** (Environment state). *We define the environment state as the conditional expectation under the environment dynamics given a realization of a state-action pair*

$$\check{S}_{t+1} := \mathbb{E}_{P_S} [S_{t+1} | S_t = s_t, A_t = a_t, W_t]. \quad (5)$$

Thus,  $\check{S}_{t+1}$  is an RV, where the randomness is induced by the process noise and has an aleatoric nature. If we additionally condition on the realization  $W_t = w_t$ , we obtain a deterministic object.

**Definition 2** (Model state). *We define the model state as the conditional expectation under  $\hat{P}_{S,TS}$*

$$\hat{S}_{t+1} := \mathbb{E}_{\hat{P}_{S,TS}} [S_{t+1} | S_t = s_t, A_t = a_t, W_t, \Theta_t]. \quad (6)$$

As discussed in Section 3, stochasticity in  $\hat{S}_{t+1}$  is induced not only by  $W_t$  but also by the randomness in the parameters  $\Theta_t$ . We project the uncertainty in the parameter space  $\vartheta$  to  $\mathcal{S}$ .

**Definition 3** (Model error process). *We define a model error process*

$$\Delta_t = \hat{S}_{t+1} - \check{S}_{t+1} \quad (7)$$

that, given a realization of process noise  $W_t = w_t$ , projects uncertainty in  $\vartheta$  to  $\mathcal{S}$

$$\mathbb{E} [\Delta_t | W_t = w_t] = \mathbb{E}_{\hat{P}_{S,TS}} [S_{t+1} | s_t, a_t, w_t, \Theta_t] - \mathbb{E}_{P_S} [S_{t+1} | s_t, a_t, w_t]. \quad (8)$$

We refer to the projected parameter uncertainty as epistemic uncertainty.

Further, we restrict model usage to a sufficiently accurate subset  $\mathcal{E} \subseteq \mathcal{S} \times \mathcal{A}$ , as proposed in Frauenknecht et al. (2024). We define  $\mathcal{E}$  amenable to the Infoprop setting in Section 4.4 and make the following assumptions when performing model-based rollouts in  $\mathcal{E}$ :

**Assumption 1** (Consistent estimator of aleatoric uncertainty). *The model's predictive variance  $\hat{\Sigma}_{\Theta_t}$  is a consistent estimator of  $\Sigma$  following the definition of Julier & Uhlmann (2001), i.e.*

$$\left( \hat{\Sigma}_{\Theta_t}(S_t, A_t) - \Sigma(S_t, A_t) \right) \succcurlyeq 0 \quad \forall (S_t, A_t) \in \mathcal{E}. \quad (9)$$

**Assumption 2** (Unbiased estimator). *The model bias  $\mu^\Delta$  is negligible. Thus  $\hat{S}_{t+1}$  according to (4) is an unbiased estimator of  $\check{S}_{t+1}$  according to (3), i.e.*

$$\mathbb{E} [\hat{S}_{t+1} | S_t, A_t] = \mathbb{E} [\check{S}_{t+1} | S_t, A_t] \quad \forall (S_t, A_t) \in \mathcal{E}. \quad (10)$$

Figure 1 empirically shows that these assumptions are reasonable. The Infoprop distribution is slightly more stochastic than the ground truth process, which indicates that Assumption 1 holds. Assumption (9) states, the model does not underestimate aleatoric uncertainty; the Infoprop rollouts should be at least as stochastic as the true process. Further, we observe no substantial bias of the Infoprop distribution underscoring the soundness of Assumption 2. Infoprop shows a similar behavior in high dimensional problems as reported in Section 6.

## 4.2 DECOMPOSING THE MODEL STATE IN SIGNAL AND NOISE

We aim to isolate the stochasticity due to parameter uncertainty in  $\hat{S}_{t+1}$ . We use the model error process (8) to project the noise in  $\vartheta$  to the same space as the signal, i.e. the dynamics, which is  $\mathcal{S}$ . The parameter distribution  $\Theta_t \sim \mathbb{P}_\Theta$  can induce arbitrarily complex distributions  $\Delta_t \sim \mathbb{P}_\Delta$ . To simplify the analysis, we solely consider the first two moments of  $\mathbb{P}_\Delta$ , namely  $\mu^\Delta$  and  $\Sigma^\Delta$ . This allows to reformulate the propagation equation (4) of the model state

$$\hat{S}_{t+1} = \check{S}_{t+1} + \Delta_t \approx \check{S}_{t+1} + \mu^\Delta(S_t, A_t) + L^\Delta(S_t, A_t)N_t \quad (11)$$

concerning the  $\check{S}_{t+1}$  and  $\Delta_t$  represented by  $\mu^\Delta(S_t, A_t)$  the model bias,  $\Sigma^\Delta(S_t, A_t)$  the epistemic variance with Cholesky decomposition  $L^\Delta(S_t, A_t)$ , and  $N_t$  the epistemic noise.

By Assumption 2, we have  $\mu^\Delta(S_t, A_t) = 0 \quad \forall (S_t, A_t) \in \mathcal{E}$ . Consequently, we can interpret the model rollout as communication through a Gaussian noise channel Cover & Thomas (2006) via (11).

Based on the propagation equation (11), we aim to infer the maximum likelihood estimate of  $\check{S}_{t+1}$  from  $E$  realizations of  $\{\mathbb{E}[\hat{S}_{t+1}|N_t = n_t^e]\}_{e=1}^E$ , to use it as the predictive distribution for our rollout scheme. As we cannot sample  $N_t$  directly, we instead use an equivalent definition of  $\hat{S}_{t+1}$ .

**Definition 4** (Model state concerning epistemic uncertainty). *Based on the model error process (8) the model state is defined as*

$$\hat{S}_{t+1} = \mathbb{E}_{\hat{P}_{\mathcal{S}, \text{TS}}} [S_{t+1}|S_t = s_t, A_t = a_t, W_t, \Delta_t] \approx \mathbb{E}_{\hat{P}_{\mathcal{S}, \text{TS}}} [S_{t+1}|S_t = s_t, A_t = a_t, W_t, N_t] \quad (12)$$

Reformulating (6) concerning  $\Delta_t$  does not change the information content or the induced sigma-algebra, as  $\Delta_t$  is a measurable function of  $\Theta_t$ . In the simplified setting of solely considering the first two moments of  $\mathbb{P}_\Delta$ ,  $N_t$  fully describes stochasticity due to model error. In reverse, we can obtain realizations  $\{\mathbb{E}[\hat{S}_{t+1}|\Theta_t = \theta_t^e]\}_{e=1}^E$  and interpret them as samples  $\{\mathbb{E}[\hat{S}_{t+1}|N_t = n_t^e]\}_{e=1}^E$ .

**Lemma 1.** *Given  $E$  realizations of  $\mathbb{E}[\hat{S}_{t+1}|\Theta_t = \theta_t^e]$ , we can estimate the environment state using maximum likelihood as*

$$\check{S}_{t+1} = \mathbb{E}[\hat{S}_{t+1}|N_t = 0] \approx \bar{S}_{t+1} = \bar{\mu}(S_t, A_t) + \bar{L}(S_t, A_t)W_t \quad (13)$$

*Proof.* see Appendix D.2.1 □

**Lemma 2.** *Following this line of thought, the maximum likelihood estimate of  $\Sigma^\Delta$  is given by*

$$\bar{\Sigma}^\Delta(S_t, A_t) = \frac{1}{E} \sum_{e=1}^E (\hat{\mu}_{\Theta_t = \theta_t^e}(S_t, A_t) - \bar{\mu}(S_t, A_t)) (\hat{\mu}_{\Theta_t = \theta_t^e}(S_t, A_t) - \bar{\mu}(S_t, A_t))^\top. \quad (14)$$

*Proof.* see Appendix D.2.2 □

Given the maximum likelihood estimates of the environment state  $\bar{S}_{t+1}$  and the epistemic variance  $\bar{\Sigma}^\Delta$ , we can decompose the model state  $\hat{S}_{t+1}$  in a signal and noise fraction according to (11) in  $\mathcal{E}$ .

## 4.3 CONSTRUCTING THE INFOPROP STATE

Having decomposed  $\hat{S}_{t+1}$  into signal  $\bar{S}_{t+1}$  and noise  $\bar{\Sigma}^\Delta$ , allows us to define the Infoprop state.

**Definition 5** (Infoprop state). *We define the Infoprop state*

$$\tilde{S}_{t+1} := \mathbb{E}[\bar{S}_{t+1}|\hat{S}_{t+1} = \hat{s}_{t+1}] = \mathbb{E}_{\hat{P}_{\mathcal{S}, \text{IP}}} [S_{t+1}|S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}, U_t] \quad (15)$$

*as the conditional expectation of the estimated environment state given a sample of the model state. We derive the corresponding Infoprop kernel  $\hat{P}_{\mathcal{S}, \text{IP}}(\cdot|S_t, A_t, \hat{S}_{t+1}) = \mathcal{N}(\bar{\mu}(S_t, A_t, \hat{S}_{t+1}), \bar{\Sigma}(S_t, A_t, \hat{S}_{t+1}))$  with the conditional noise  $U_t \sim \mathcal{N}(0, I)$  in Appendix D.3.*

Consequently, the Infoprop state aims to infer the signal  $\bar{S}_{t+1}$  given a noisy observation  $\hat{s}_{t+1}$ . Propagating model-based rollouts using  $\tilde{S}_{t+1}$ , yields favorable properties as stated in Theorem 1.

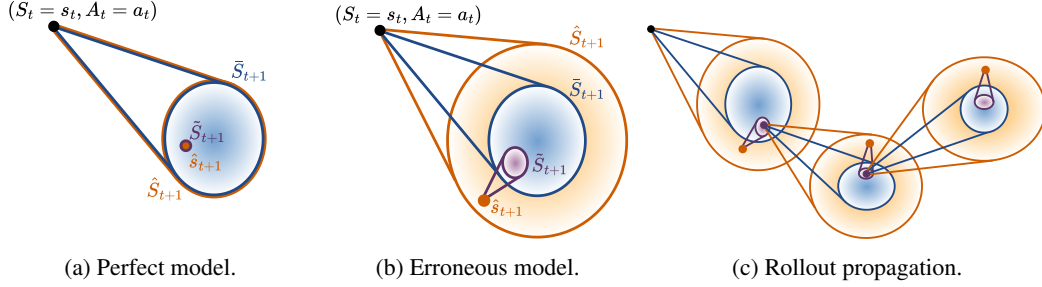


Figure 3: Infoprop rollout mechanism. (a), (b): Generating the Infoprop state  $\tilde{\tilde{S}}_{t+1}$  from the estimated predictive distribution  $\tilde{S}_{t+1}$  and the model sample  $\hat{s}_{t+1}$ . (c) Performing an Infoprop rollout.

**Theorem 1** (Infoprop state). *By construction,  $\tilde{\tilde{S}}_{t+1}$  addresses questions (i) and (ii) of Section 3.*

(i) *The distribution of Infoprop states is identical to the estimated environment distribution.*

$$\tilde{\tilde{S}}_{t+1} \stackrel{\text{dist}}{=} \tilde{S}_{t+1} \quad (16)$$

*Proof.* see Appendix D.4. □

(ii) *The sum of marginal entropies of  $\tilde{\tilde{S}}_{t+1}$  defines the information loss along an Infoprop rollout.*

$$\mathbb{H}\left(\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_T \mid S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T\right) = \sum_{t=0}^T \mathbb{H}\left(\tilde{\tilde{S}}_{t+1}\right) \quad (17)$$

*Proof.* see Appendix D.5. □

Figure 3 illustrates the Infoprop rollout mechanism and provides intuition for Theorem 1. In the case of a perfect model, i.e.  $\bar{\Sigma}^\Delta = 0$ , depicted in Figure 3a, the realization  $\hat{s}_{t+1}$  provides the information about the process noise realization  $w_t$  without ambiguity. Consequently, the belief about the environment state given the sample from the model  $\tilde{S}_{t+1} = \mathbb{E}[\tilde{S}_{t+1} \mid \hat{S}_{t+1} = \hat{s}_{t+1}]$  is a deterministic object and  $\mathbb{H}(\tilde{S}_{t+1}) = 0$ . In the general scenario of  $\bar{\Sigma}^\Delta > 0$  depicted in Figure 3b, the epistemic uncertainty results in ambiguity about the environment state given  $\hat{s}_{t+1}$ , such that  $\mathbb{H}(\tilde{S}_{t+1}) > 0$ . Notably, conditioning  $\tilde{S}_{t+1}$  on  $\hat{s}_{t+1}$ , results in Infoprop predictions  $\tilde{\tilde{S}}_{t+1}$  following estimated environment distribution  $\tilde{S}_{t+1}$  as stated in Theorem 1 (i). This results in a data distribution that closely resembles the environment dynamics as desired in question (i) of Section 3. Finally, Figure 3c depicts a Infoprop rollout propagated via realization  $\tilde{s}_{t+1}$ . We measure data corruption due to model error using the conditional entropy of a rollout under the estimated environment dynamics  $(\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_T)$  given the realizations observed from the model  $(S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T)$ , i.e. *given the observed model trajectory, how sure are we on how the corresponding environment trajectory would look like?*. As per Theorem 1 (ii), this trajectory-based approach to uncertainty can be addressed with the accumulated marginal entropy of  $\tilde{\tilde{S}}_{t+1}$ , addressing question (ii) of Section 3.

#### 4.4 ROLLOUT TERMINATION CRITERIA

Having introduced how to propagate Infoprop rollouts, the question remains when to terminate them. In the following, we propose two termination criteria to address question (iii) of Section 3.

First, Infoprop rollouts build on the assumption that model usage is restricted to a sufficiently accurate subset  $\mathcal{E} \subseteq \mathcal{S} \times \mathcal{A}$ , following the ideas of Frauenknecht et al. (2024).

**Definition 6** (Sufficiently accurate subset). *We define the sufficiently accurate subset*

$$\mathcal{E} := \{(s_t, a_t) \in \mathcal{S} \times \mathcal{A} \mid \mathbb{H}(\tilde{\tilde{S}}_{t+1}) \leq \lambda_1, \hat{s}_{t+1} \sim \hat{P}_{S, \text{TS}}(\cdot \mid s_t, a_t)\} \quad (18)$$

*based on a threshold  $\lambda_1$  for the single-step information loss  $\mathbb{H}(\tilde{\tilde{S}}_{t+1})$ .*

Second, we restrict Infoprop rollouts to sufficiently accurate paths to limit uncertainty accumulation.

**Definition 7** (Sufficiently accurate path). *Based on the estimated information loss along a rollout (17), we define the set of sufficiently accurate paths of length  $t' \in \{1, \dots, T\}$  as*

$$\mathcal{P}^{t'} := \left\{ (s_t, a_t)_{t=0}^{t'} \in (\mathcal{S} \times \mathcal{A})^{t'} \left| \sum_{t=0}^{t'} \mathbb{H}(\tilde{S}_{t+1}) \leq \lambda_2 \right. \right\}. \quad (19)$$

Heuristics for determining values of  $\lambda_1$  and  $\lambda_2$  depend on the class of AES model and MBRL algorithm at hand with an example provided in Section 5. Combining the steps above yields the Infoprop rollout mechanism illustrated in Algorithm 1.

---

**Algorithm 1** Infoprop

---

**Require:**  $s_0$   
**while**  $t < T + 1$  **do**  
     $a_t \sim \pi(\cdot | s_t)$   
    **for**  $e \in \{1, \dots, E\}$  **do**  
         $\theta_t^e \sim \mathbb{P}_\Theta$   
         $\bar{S}_{t+1}(s_t, a_t)$  from (13), and  $\bar{\Sigma}^\Delta(s_t, a_t)$  from (14)  
         $\hat{s}_{t+1} = \mathbb{E} \left[ \hat{S}_{t+1} | W_t = w_t, \Theta_t = \theta_t^e \right]$  with  $w_t \sim \mathcal{N}(0, I)$ ,  $\theta_t^e \sim \mathcal{U}(\{\theta_t^1, \dots, \theta_t^E\})$   
         $\tilde{S}_{t+1}$  from (51) and  $\mathbb{H}(\tilde{S}_{t+1})$  from (24)  
        **if**  $\mathbb{H}(\tilde{S}_{t+1}) > \lambda_1$  **then**  
            **break**  
        **else if**  $\sum_{t'=0}^t \mathbb{H}(\tilde{S}_{t'+1}) > \lambda_2$  **then**  
            **break**  
        **else**  
             $s_t \leftarrow \mathbb{E}[\tilde{S}_{t+1} | U_t = u_t]$  with  $u_t \sim \mathcal{N}(0, I)$

---

## 5 AUGMENTING STATE-OF-THE-ART: INFOPROP-DYNA

While the Infoprop rollout mechanism is applicable to different kinds of MBRL with AES models, we illustrate its capabilities in a Dyna-style architecture with probabilistic ensemble (PE) models Lakshminarayanan et al. (2017). We design *Infoprop-Dyna* by integrating the Infoprop rollout mechanism in the state-of-the-art framework proposed in Janner et al. (2019) with minor adaptations.

As discussed in Section 4.4, heuristics for  $\lambda_1$  and  $\lambda_2$  depend on the algorithm at hand. In Infoprop-Dyna, we take the common approach Chua et al. (2018); Janner et al. (2019) of neglecting cross-correlations between state dimensions for computational reasons. Thus, we can consider data corruption of each state dimension independently. As the predictive quality of different state dimensions can differ substantially, we choose both thresholds as  $n_{\mathcal{S}}$  dimensional vectors, such that a rollout is terminated as soon as the data corruption of any dimension overshoots the corresponding threshold.

In Dyna-style MBRL Janner et al. (2019), the dynamics model is trained on the data distribution observed during environment interaction. The corresponding transitions are stored in an environment replay buffer  $\mathcal{D}_{\text{env}} = \{(\check{s}_t^{(b)}, \check{a}_t^{(b)}, \check{r}_{t+1}^{(b)}, \check{s}_{t+1}^{(b)})\}_{b=1}^{|\mathcal{D}_{\text{env}}|}$ , where  $(b)$  indicates the index in the replay buffer. After a fixed number of interaction steps between a model-free RL agent and the environment, the dynamics model is retrained on the data in  $\mathcal{D}_{\text{env}}$ , model-based rollouts are performed, and the data is stored in a replay buffer  $\mathcal{D}_{\text{mod}}$  to train the model-free RL agent. Consequently, we assume the PE model to be accurate within the data distribution of  $\mathcal{D}_{\text{env}}$  and build the heuristic for  $\lambda_1$  and  $\lambda_2$  on the predictive uncertainty within the environment buffer.

After each round of retraining the PE model, we compute a set of dimension-wise Infoprop state entropies for single-step predictions in  $\mathcal{D}_{\text{env}}$  according to

$$\mathcal{H}^k = \left\{ \mathbb{H} \left( \bar{S}_{t+1}^k | S_t = \check{s}_t^{(b)}, A_t = \check{a}_t^{(b)}, \hat{S}_{t+1}^k = \hat{s}_{t+1}^{k,(b)} \right) = \mathbb{H} \left( \tilde{S}_{t+1}^{k,(b)} \right) \right\}_{b=1}^{|\mathcal{D}_{\text{env}}|} \quad (20)$$

where  $k \in \{1, \dots, n_S\}$  indicates the corresponding state dimension. We define the dimension-wise thresholds  $\lambda_1^k$  and  $\lambda_2^k$  based on the cumulative distribution function of dimension-wise entropies

$$F_{\mathcal{H}^k}(h) = \frac{1}{|\mathcal{H}^k|} \sum_{h' \in \mathcal{H}^k} \mathbb{1}[h' \leq h]. \quad (21)$$

The  $k^{\text{th}}$  element of  $\lambda_1$  is defined as the  $\zeta_1$  quantile of the single-step entropy set

$$\lambda_1^k = \inf \{h \in \mathcal{H}^k : F_{\mathcal{H}^k}(h) \geq \zeta_1\} \quad (22)$$

and limits model usage to the sufficiently accurate subset  $\mathcal{E}$ . To restrict rollouts of length  $t'$  to  $\mathcal{P}^{t'}$ , we define the  $k^{\text{th}}$  element of  $\lambda_2$  as the  $\zeta_2$  quantile of the entropy set scaled by  $\xi$

$$\lambda_2^k = \xi \inf \{h \in \mathcal{H}^k : F_{\mathcal{H}^k}(h) \geq \zeta_2\}. \quad (23)$$

Here,  $\zeta_2$  denotes a quantile corresponding to precise predictions and  $\xi$  to the number of prediction steps we are willing to accumulate the resulting data corruption. We choose  $\zeta_1 = 0.99$ ,  $\zeta_2 = 0.01$  and  $\xi = 100$  for all experiments in Section 6 without further hyperparameter tuning.

We use pink noise for environment exploration Eberhard et al. (2023) to quickly expand  $\mathcal{E}$  Frauenknecht et al. (2024). Pseudocode is provided in Algorithm 3 of Appendix C.

## 6 EXPERIMENTS AND DISCUSSION

To demonstrate the benefits of the Infoprop mechanism, we compare Infoprop-Dyna to state-of-the-art Dyna-style MBRL algorithms on MuJoCo Todorov et al. (2012) benchmark tasks. We report

- substantial improvements in the consistency of predicted data, especially over long horizons;
- effective rollout termination based on accumulated model error propagation; and
- state-of-the-art performance in Dyna-style MBRL on several MuJoCo tasks.

Furthermore, we discuss the limitations of naively integrating Infoprop into the standard Dyna-style setup Janner et al. (2019) and point to further research questions.

### 6.1 EXPERIMENTAL SETUP

We compare Infoprop-Dyna to Model-Based Policy Optimization (MBPO) Janner et al. (2019) and Model-Based Actor-Critic with Uncertainty-Aware Rollout Adaption (MACURA) Frauenknecht et al. (2024) as well as to Soft Actor-Critic (SAC) Haarnoja et al. (2018) that represents the model-free learner of all the Dyna-style approaches above. We build our implementation<sup>2</sup> on the code base<sup>3</sup> provided by Frauenknecht et al. (2024). Further details are provided in Appendix E.1

### 6.2 PREDICTION QUALITY

To compare different rollout mechanisms, we train an Infoprop-Dyna agent on hopper for 120000 environment interactions and perform model rollouts from states in  $\mathcal{D}_{\text{env}}$ .

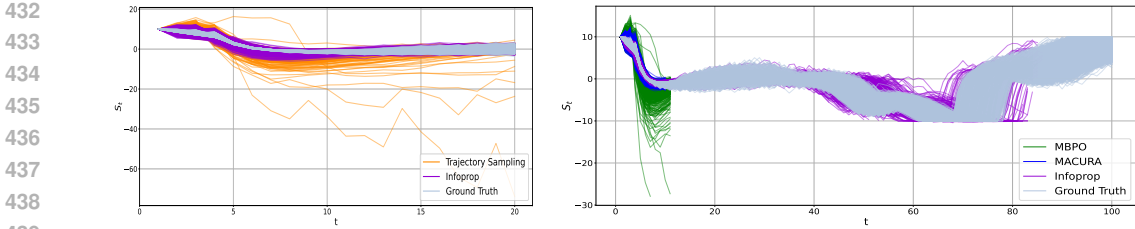
First, we evaluate the consistency of Infoprop and TS rollouts, propagating 20 steps without termination. Figure 4a depicts the resulting distributions for the 11<sup>th</sup> dimension of the hopper state. Infoprop rollouts show substantially improved data consistency compared to TS rollouts, underscoring the ability of Infoprop to effectively mitigate model error propagation.

Next, we compare the rollout mechanisms of MBPO and MACURA based on TS sampling with Infoprop-Dyna rollouts. Figure 4b shows the results for 11<sup>th</sup> dimension of the hopper and a maximum rollout length of 100 steps. MBPO rollouts are propagated for 11 steps following the schedule proposed in Janner et al. (2019), resulting in a widely spread distribution. In contrast, MACURA has an adaptive rollout length capped at 10 steps Frauenknecht et al. (2024), leading to better data consistency. The improved predictive distribution and capability to estimate accumulated error of Infoprop allows for substantially longer rollouts up to 100 steps. The Infoprop termination criteria reliably stop distorted rollouts, resulting in consistent rollouts over long horizons. Appendix E.2 provides additional results for setting the maximum rollout length of all three approaches to 100.

<sup>2</sup>Code will be published upon acceptance and is currently provided in the supplementary material.

<sup>3</sup><https://github.com/Data-Science-in-Mechanical-Engineering/macura>

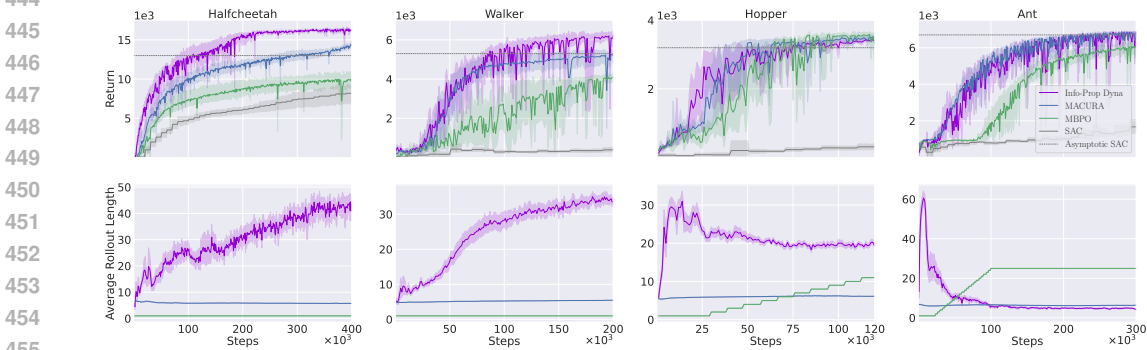




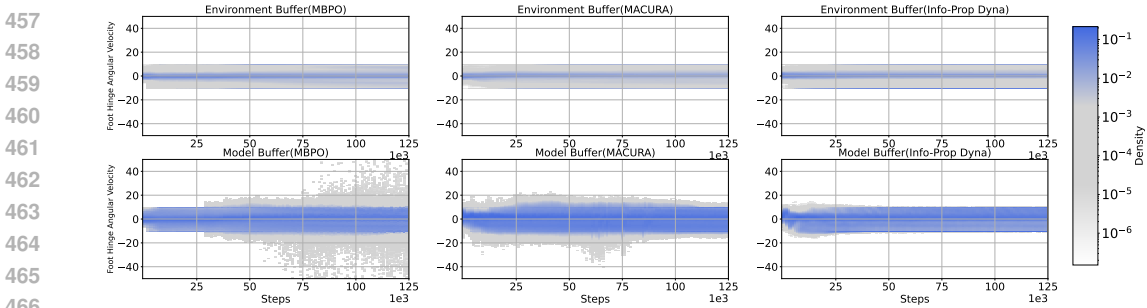
(a) Trajectory Sampling vs. Infoprop

(b) MBPO vs. MACURA vs. Infoprop-Dyna

Figure 4: Predictive quality of rollouts in the 11<sup>th</sup> state dimension of MuJoCo hopper. (a) Rollouts according to Trajectory Sampling (TS) and Infoprop. (b) Rollout schemes of MBPO and MACURA based on TS compared to Infoprop-Dyna.



(a) Performance and average rollout length on MuJoCo tasks.



(b) Adequacy of  $\mathcal{D}_{\text{mod}}$  on the 11<sup>th</sup> state dimension of hopper.

Figure 5: Evaluation on MuJoCo tasks. (a) Infoprop-Dyna shows state-of-the-art performance for Dyna-style MBRL on several MuJoCo tasks while considerably increasing average rollout length on most tasks. (b) Infoprop-Dyna shows substantially improved consistency between  $\mathcal{D}_{\text{env}}$  and  $\mathcal{D}_{\text{mod}}$ .

### 6.3 PERFORMANCE EVALUATION

As depicted in the top row of Figure 5a, Infoprop-Dyna performs on par with or better than MACURA, while substantially outperforming MBPO with respect to data efficiency and asymptotic performance. Notably, Infoprop-Dyna consistently outperforms SAC with a fraction of environment interaction. The bottom row of Figure 5a depicts the average rollout lengths. Infoprop-Dyna shows substantially increased rollout lengths compared to prior methods in all environments but ant.

A major concern of this work is the consistency of model-based rollouts with the environment distribution. Figure 5b depicts the data distribution in  $\mathcal{D}_{\text{env}}$  and  $\mathcal{D}_{\text{mod}}$  of the respective Dyna-style approaches throughout training for the 11<sup>th</sup> dimension of the hopper state. The figure shows a histogram over state values over the course of training. It can be seen that the model data distribution of Infoprop-Dyna closely follows the distribution observed in the environment, while both the data from MBPO and MACURA show severe outliers. This is the case, even though the rollout data in Infoprop-Dyna is obtained from substantially longer rollouts as can be seen from Figure 5a which indicates the capabilities of the Infoprop rollout mechanism.

## 6.4 LIMITATIONS AND OUTLOOK

Despite the excellent quality of model-generated data with the Infoprop rollout, the limitations of Infoprop-Dyna are most apparent on MuJoCo humanoid with results provided in E.3. These show instabilities in learning and point to structural problems when integrating Infoprop rollouts naively into standard Dyna-style architectures Janner et al. (2019).

Figure 5b shows that the long rollouts of Infoprop-Dyna can cause rapid distribution shifts in  $\mathcal{D}_{\text{mod}}$ , especially early in training. These nonstationary buffers are a challenge to deep Q-learning methods Mnih et al. (2015). Another issue is primacy bias in model learning Qiao et al. (2023), where the model overfits to initial data and subsequently struggles to generalize, as seen in the decreasing rollout length for the ant environment in Figure 5a. The main problem with Infoprop-Dyna is likely overfitting critics and plasticity loss Nikishin et al. (2022); D’Oro et al. (2023), as also reported by Frauenknecht et al. (2024) for Dyna-style MBRL trained on high-quality data. We provide an ablation on this observation and sketch methods to counteract this phenomenon in Appendix E.4.

## 7 RELATED WORK

The negative effects of accumulated model error on the performance of MBRL methods is a long-studied problem Venkatraman et al. (2015); Talvitie (2016); Asadi et al. (2018b;a).

Different model architectures have been proposed to mitigate this issue, such as trajectory models Asadi et al. (2019); Lambert et al. (2021), bidirectional models Lai et al. (2020), temporal segment models Mishra et al. (2017) or self-correcting models Talvitie (2016). These architectures, however, imply substantial additional effort for model learning, such that state-of-the-art performance in the respective fields of MBRL is often reported for simpler single-step model architectures Chua et al. (2018); Janner et al. (2019); Buckman et al. (2018).

These approaches address the problem of error accumulation by keeping model-based rollouts sufficiently short. Janner et al. (2019) introduce the concept of branched rollouts that allows to cover relevant parts of  $\mathcal{S}$  with short model rollouts. Other methods weight rollouts of different lengths according to their single-step uncertainty Buckman et al. (2018) or use single-step uncertainty to schedule rollout length Pan et al. (2020); Frauenknecht et al. (2024). Infoprop allows to infer model data consistent with the environment distribution over long rollout horizons using comparatively simple model architectures and computationally cheap conditioning operations.

Infoprop is inspired by an information-theoretic view on RL Lu et al. (2023). Thus far, information-theoretic arguments have been mostly used to improve the exploration Haarnoja et al. (2018); Lu & Roy (2019); Ahmed et al. (2019); Mohamed & Rezende (2015) and generalization Tishby & Zaslavsky (2015); Lu et al. (2020); Igl et al. (2019); Islam et al. (2023) of model-free RL methods. While aspects of dynamical systems such as causality, modeling, and control Lozano-Duran & Aranz (2021), predictability Kleeman (2011) or dealing with noisy observations Gattami (2014) have been studied from an information theoretic perspective, these works do not directly apply to the MBRL setup nor extend to long model-based rollouts.

## 8 CONCLUDING REMARKS

Data consistency of model-based rollouts is a key criterion for the performance of MBRL approaches. This work proposes the novel Infoprop mechanism that substantially improves rollouts with common AES models. We reduce the influence of epistemic uncertainty on the predictive distribution of model-based rollouts, keep track of data corruption through propagated model error over long horizons, and terminate rollouts based on data corruption. This allows for considerably increased rollout lengths while substantially improving data consistency simultaneously.

While Infoprop is applicable to a broad range of MBRL methods, we demonstrate its capabilities by naively integrating Infoprop into a standard Dyna-style MBRL architecture Janner et al. (2019) resulting in the Infoprop-Dyna algorithm. We report state-of-the-art performance in several MuJoCo tasks while pointing to necessary adaptations to the existing algorithmic framework to fully unleash the potential of Infoprop rollouts.

## 540 REFERENCES

- 541  
542 Zafarali Ahmed, Nicolas Le Roux, Mohammad Norouzi, and Dale Schuurmans. Understanding the  
543 impact of entropy on policy optimization. *Int. Conf. on Machine Learning*, abs/1811.11214, 2019.
- 544 Kavosh Asadi, Evan Cater, Dipendra Misra, and Michael L. Littman. Towards a Simple Approach  
545 to Multi-step Model-based Reinforcement Learning. *arXiv*, October 2018a. doi: 10.48550/arXiv.  
546 1811.00128.
- 547 Kavosh Asadi, Dipendra Misra, and Michael L. Littman. Lipschitz Continuity in Model-based  
548 Reinforcement Learning. *arXiv*, April 2018b. doi: 10.48550/arXiv.1804.07193.
- 549 Kavosh Asadi, Dipendra Misra, Seungchan Kim, and Michel L. Littman. Combating the  
550 Compounding-Error Problem with a Multi-step Model. *arXiv*, May 2019. doi: 10.48550/arXiv.  
551 1905.13320.
- 552 Philipp Becker and Gerhard Neumann. On Uncertainty in Deep State Space Models for Model-  
553 Based Reinforcement Learning. *Transactions on Machine Learning Research*, October 2022.  
554 doi: 10.48550/arXiv.2210.09256.
- 555 Marc G. Bellemare, Will Dabney, and Mark Rowland. *Distributional Reinforcement Learning*. MIT  
556 Press, 2023.
- 557 Thomas Bi and Raffaello D’Andrea. Sample-efficient learning to solve a real-world labyrinth game  
558 using data-augmented model-based reinforcement learning. In *IEEE Int. Conf. on Robotics and  
559 Automation*, pp. 7455–7460, 2024. doi: 10.1109/ICRA57147.2024.10610577.
- 560 Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-  
561 efficient reinforcement learning with stochastic ensemble value expansion. In *Int. Conf. on Neural  
562 Information Processing Systems*. 2018.
- 563 Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep Reinforcement  
564 Learning in a Handful of Trials using Probabilistic Dynamics Models. *Adv. in Neural Information  
565 Processing Systems*, 2018.
- 566 Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory 2nd Edition (Wiley Series in  
567 Telecommunications and Signal Processing) by Thomas M. Cover Joy A. Thomas(2006-07-18)*.  
568 Wiley-Interscience, Hoboken, NJ, USA, January 2006.
- 569 Marc Peter Deisenroth and Carl Edward Rasmussen. PILCO: a model-based and data-efficient  
570 approach to policy search. In *Int. Conf. on Machine Learning*. 2011.
- 571 Pierluca D’Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G. Bellemare, and  
572 Aaron C. Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier.  
573 In *Int. Conf. on Learning Representations*, 2023.
- 574 Onno Eberhard, Jakob Hollenstein, Cristina Pinneri, and Georg Martius. Pink noise is all you need:  
575 Colored noise exploration in deep reinforcement learning. In *Int. Conf. on Learning Representa-  
576 tions*, 2023.
- 577 Vladimir Feinberg, Alvin Wan, Ion Stoica, Michael I. Jordan, Joseph E. Gonzalez, and Sergey  
578 Levine. Model-Based Value Estimation for Efficient Model-Free Reinforcement Learning. *Int.  
579 Conf. on Machine Learning*, 2018.
- 580 Bernd Frauenknecht, Artur Eisele, Devdutt Subhasish, Friedrich Solowjow, and Sebastian Trimpe.  
581 Trust the Model Where It Trusts Itself – Model-Based Actor-Critic with Uncertainty-Aware Roll-  
582 out Adaption. *Int. Conf. on Machine Learning*, May 2024. doi: 10.48550/arXiv.2405.19014.
- 583 Ather Gattami. Kalman meets shannon. *ArXiv*, 2014.
- 584 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash  
585 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic Algo-  
586 rithms and Applications. *arXiv*, 2018.
- 587  
588  
589  
590  
591  
592  
593

- 594 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James  
595 Davidson. Learning Latent Dynamics for Planning from Pixels. In *Int. Conf. on Machine Learn-*  
596 *ing*. PMLR, May 2019.
- 597 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to Control: Learn-  
598 ing Behaviors by Latent Imagination. In *Int. Conf. on Learning Representations*, April 2020.
- 600 Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari with  
601 Discrete World Models. *Int. Conf. on Learning Representations*, October 2021.
- 602 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering Diverse Domains  
603 through World Models. *Int. Conf. on Learning Representations*, January 2023.
- 605 Maximilian Igl, Kamil Ciosek, Yingzhen Li, Sebastian Tschitschek, Cheng Zhang, Sam Devlin,  
606 and Katja Hofmann. Generalization in reinforcement learning with selective noise injection and  
607 information bottleneck. In *Int. Conf. on Neural Information Processing Systems*, 2019.
- 608 Riashat Islam, Hongyu Zang, Manan Tomar, Aniket Didolkar, Md Mofijul Islam, Samin Yeasar  
609 Arnob, Tariq Iqbal, Xin Li, Anirudh Goyal, Nicolas Heess, and Alex Lamb. Representation  
610 learning in deep rl via discrete information bottleneck. In *Int. Conf. on Artificial Intelligence and*  
611 *Statistics*, 2023.
- 613 Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: model-  
614 based policy optimization. In *Int. Conf. on Neural Information Processing Systems*. 2019.
- 615 Simon Julier and Jeffrey Uhlmann. *General Decentralized Data Fusion with Covariance Intersec-*  
616 *tion (CI)*. 06 2001. doi: 10.1201/9781420038545.ch12.
- 618 Elia Kaufmann, Leonard Bauersfeld, Antonio Loquercio, Matthias Müller, Vladlen Koltun, and  
619 Davide Scaramuzza. Champion-level drone racing using deep reinforcement learning. *Nature*,  
620 August 2023.
- 621 Richard Kleeman. Information theory and dynamical system predictability. *Entropy*, 2011. doi:  
622 10.3390/e13030612.
- 623 Hang Lai, Jian Shen, Weinan Zhang, and Yong Yu. Bidirectional Model-based Policy Optimization.  
624 In *Int. Conf. on Machine Learning*. PMLR, 2020.
- 626 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
627 uncertainty estimation using deep ensembles. In *Int. Conf. on Neural Information Processing*  
628 *Systems*. 2017.
- 629 Nathan O. Lambert, Albert Wilcox, Howard Zhang, Kristofer S. J. Pister, and Roberto Calandra.  
630 Learning Accurate Long-term Dynamics for Model-based Reinforcement Learning. *IEEE Conf*  
631 *on Decision and Control*, December 2021. doi: 10.48550/arXiv.2012.09156.
- 633 Adrian Lozano-Duran and Gonzalo Arranz. Information-theoretic formulation of dynamical sys-  
634 tems: causality, modeling, and control. *ArXiv*, 2021.
- 635 Xingyu Lu, Kimin Lee, P. Abbeel, and Stas Tiomkin. Dynamics generalization via information  
636 bottleneck in deep reinforcement learning. *ArXiv*, 2020.
- 637  
638 Xiuyuan Lu and Benjamin Van Roy. Information-theoretic confidence bounds for reinforcement  
639 learning. *Int. Conf. on Neural Information Processing Systems*, 2019.
- 640 Xiuyuan Lu, Benjamin Van Roy, Vikranth Dwaracherla, Morteza Ibrahimi, Ian Osband, and Zheng  
641 Wen. Reinforcement learning, bit by bit. *Foundations and Trends® in Machine Learning*, 2023.  
642 doi: 10.1561/22000000097.
- 643  
644 Carlos E. Luis, Alessandro G. Bottero, Julia Vinogradska, Felix Berkenkamp, and Jan Peters.  
645 Model-Based Epistemic Variance of Values for Risk-Aware Policy Optimization. *arXiv*, 2023.
- 646  
647 Nikhil Mishra, Pieter Abbeel, and Igor Mordatch. Prediction and Control with Temporal Segment  
Models. *Int. Conf. on Machine Learning*, March 2017. doi: 10.48550/arXiv.1703.04070.

- 648 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-  
649 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,  
650 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan Wier-  
651 stra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.  
652 *Nature*, 2015.
- 653 Shakir Mohamed and Danilo J. Rezende. Variational information maximisation for intrinsically  
654 motivated reinforcement learning. In *Int. Conf. on Neural Information Processing Systems*, 2015.  
655
- 656 Anusha Nagabandi, Gregory Kahn, Ronald S. Fearing, and Sergey Levine. Neural Network Dynam-  
657 ics for Model-Based Deep Reinforcement Learning with Model-Free Fine-Tuning. In *IEEE Int.*  
658 *Conf. on Robotics and Automation*. 2018.
- 659 Michal Nauman, Michał Borkiewicz, Piotr Miłoś, Tomasz Trzcinski, Mateusz Ostaszewski, and  
660 Marek Cygan. Overestimation, Overfitting, and Plasticity in Actor-Critic: the Bitter Lesson of  
661 Reinforcement Learning. In *International Conference on Machine Learning*, pp. 37342–37364.  
662 PMLR, July 2024. URL [https://proceedings.mlr.press/v235/nauman24a.](https://proceedings.mlr.press/v235/nauman24a.html)  
663 [html](https://proceedings.mlr.press/v235/nauman24a.html).  
664
- 665 Evgenii Nikishin, Max Schwarzer, Pierluca D’Oro, Pierre-Luc Bacon, and Aaron C. Courville. The  
666 primacy bias in deep reinforcement learning. In *Int. Conf. on Machine Learning*, 2022.
- 667 OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw De-  
668 biak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal  
669 Jozefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé  
670 de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon  
671 Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep  
672 reinforcement learning. 2019.
- 673 Feiyang Pan, Jia He, Dandan Tu, and Qing He. Trust the model when it is confident: masked  
674 model-based actor-critic. In *Int. Conf. on Neural Information Processing Systems*. 2020.  
675
- 676 Zhongjian Qiao, Jiafei Lyu, and Xiu Li. Mind the Model, Not the Agent: The Primacy Bias in  
677 Model-based RL. *arXiv*, October 2023. doi: 10.48550/arXiv.2310.15017.
- 678 C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, July  
679 1948.  
680
- 681 Dan Simon. *Optimal State Estimation*. January 2006. ISBN 978-0-47170858-2. doi: 10.1002/  
682 0470045345.
- 683 Laura Smith, Ilya Kostrikov, and Sergey Levine. Demonstrating A Walk in the Park: Learning to  
684 Walk in 20 Minutes With Model-Free Reinforcement Learning. *Robotics, Science and Systems*,  
685 August 2023. doi: 10.48550/arXiv.2208.07860.  
686
- 687 Richard S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *SIGART*  
688 *Bull.*, 1991.
- 689 Erik Talvitie. Self-Correcting Models for Model-Based Reinforcement Learning. *arXiv*, December  
690 2016. doi: 10.48550/arXiv.1612.06018.  
691
- 692 Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. *IEEE*  
693 *Information Theory Workshop (ITW)*, 2015.  
694
- 695 Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based control.  
696 In *Int. Conf. on Intelligent Robots and Systems*. IEEE, 2012.
- 697 Miguel Vasco, Takuma Seno, Kenta Kawamoto, Kaushik Subramanian, Peter R. Wurman, and Peter  
698 Stone. A Super-human Vision-based Reinforcement Learning Agent for Autonomous Racing in  
699 Gran Turismo. *Reinforcement Learning Conference*, June 2024. doi: 10.48550/arXiv.2406.12563.  
700
- 701 Arun Venkatraman, Martial Hebert, and J. Bagnell. Improving Multi-Step Prediction of Learned  
Time Series Models. *AAAI*, February 2015. doi: 10.1609/aaai.v29i1.9590.

702 Grady Williams, Nolan Wagener, Brian Goldfain, Paul Drews, James M. Rehg, Byron Boots, and  
703 Evangelos A. Theodorou. Information theoretic MPC for model-based reinforcement learning. In  
704 *IEEE Int. Conf. on Robotics and Automation*. IEEE, 2017.

705  
706 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn,  
707 and Tengyu Ma. Mopo: Model-based offline policy optimization. In *Adv. in Neural Information*  
708 *Processing Systems*, 2020.

709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## A NOTATION

### A.1 OBJECTS

- $S_t$  Random variable of a general state
- $\tilde{S}_t$  Random variable of the environment state
- $\bar{S}_t$  Random variable of the estimated environment state
- $\hat{S}_t$  Random variable of the model state
- $A_t$  Random variable of the action

... ..

[we will finish this for a potential camera-ready version.](#)

## B TOY EXAMPLE

In Figure 1, we illustrate the data consistency of Trajectory Sampling Chua et al. (2018) and Infoprop in a one-dimensional random walk example with  $\mathcal{S} \subseteq \mathbb{R}$  and  $\mathcal{A} \subseteq \mathbb{R}$ . The dynamics follow (3) with  $\mu(S_t, A_t) = S_t + A_t$  and  $L(S_t, A_t) = 0.01$ . Actions are distributed according to  $A_t \sim \mathcal{N}(0, 0.1)$ . All rollouts start from  $s_0 = 0$  and are propagated for 100 steps. We perform 1000 rollouts under the environment dynamics and train a Probabilistic Ensemble Lakshminarayanan et al. (2017) model according to the information provided in Table 1. Subsequently, we perform 1000 model-based rollouts with this model and the respective rollout mechanism.

| Hyperparameter             | Value   |
|----------------------------|---------|
| number of ensemble members | 5       |
| number of hidden neurons   | 2       |
| number of layers           | 1       |
| learning rate              | 0.001   |
| weight decay               | 0.00001 |
| number of epochs           | 4       |

Table 1: Hyperparameters used for training the model on the random walk dataset.

## C PSEUDOCODE ALGORITHMS

---

### Algorithm 2 Trajectory Sampling Chua et al. (2018)

---

**Require:**  $s_0$   
**while**  $t < T + 1$  **do**  
 $a_t \sim \pi(\cdot | s_t)$   
 $\hat{s}_{t+1} = \mathbb{E} \left[ \hat{\mathcal{S}}_{t+1} | W_t = w_t, \Theta_t = \theta_t \right]$  with  $w_t \sim \mathcal{N}(0, I)$  and  $\theta_t \sim \mathbb{P}_\Theta$   
 $s_t \leftarrow \hat{s}_{t+1}$

---



---

### Algorithm 3 Infoprop-Dyna (Pseudocode adapted from Janner et al. (2019))

---

**Require:** Policy  $\pi$ , predictive AES model  $p_\Theta$ , environment buffer  $\mathcal{D}_{\text{env}}$ , model buffer  $\mathcal{D}_{\text{mod}}$ , rollout parameters  $T, \zeta_1, \zeta_2, \xi$   
**for**  $N$  epochs **do**  
  **for**  $J$  steps **do**  
    Interact with the environment according to  $\pi$ ; add to  $\mathcal{D}_{\text{env}}$   
    Train model  $p_\Theta$  on  $\mathcal{D}_{\text{env}}$   
    Perform single-step predictions with  $p_\Theta$  in  $\mathcal{D}_{\text{env}}$   
    Compute  $\lambda_1$  (22) and  $\lambda_2$  (23)  
    **for**  $M$  model rollouts **do**  
      Sample  $s_0$  uniformly from  $\mathcal{D}_{\text{env}}$   
      Perform Infoprop rollouts according to Algorithm 1; add to  $\mathcal{D}_{\text{mod}}$   
    **for**  $G \cdot J$  gradient updates **do**  
      Update  $\pi$  on  $\mathcal{D}_{\text{env}} \cup \mathcal{D}_{\text{mod}}$

---



## 864 D DERIVATIONS

### 865 D.1 QUANTIZED ENTROPY

866 For a RV  $Z \in \mathcal{Z} \subseteq \mathbb{R}^{n_z}$  with  $Z \sim \mathcal{N}(\mu_Z, \Sigma_Z)$  and discretization step size  $\Delta z^{(k)}$  of the  $k^{\text{th}}$   
 867 dimension, the quantized entropy Cover & Thomas (2006) is

$$868 \mathbb{H}(Z) = \frac{1}{2} \log_2 ((2\pi e)^{n_z} |\Sigma_Z|) - \sum_{k=1}^{n_z} \log_2 \left( \Delta z^{(k)} \right). \quad (24)$$

### 874 D.2 MAXIMUM LIKELIHOOD PREDICTIVE DISTRIBUTION

#### 875 D.2.1 PROOF OF LEMMA 1

876 *Proof.* We introduce the conditional expectation over the next state under the model, given a realization  $\theta_t^e$

$$877 \hat{S}_{t+1}^e := \mathbb{E}_{\hat{P}_{S,TS}} \left[ \hat{S}_{t+1} | \Theta_t = \theta_t^e \right]. \quad (25)$$

878 Further,  $\hat{\mu}^e := \hat{\mu}_{\Theta_t = \theta_t^e}$ ,  $\hat{\Sigma}^e := \hat{\Sigma}_{\Theta_t = \theta_t^e}$  and  $\hat{L}^e := \hat{L}_{\Theta_t = \theta_t^e}$  such that

$$879 \hat{S}_{t+1}^e = \hat{\mu}^e(S_t, A_t) + \hat{L}^e(S_t, A_t) W_t. \quad (26)$$

880 Given  $E$  RVs  $\hat{S}_{t+1}^e$  we define their joint distribution

$$881 \begin{pmatrix} \hat{S}_{t+1}^1 \\ \vdots \\ \hat{S}_{t+1}^E \end{pmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \hat{\mu}^1 \\ \vdots \\ \hat{\mu}^E \end{bmatrix}, \begin{bmatrix} \hat{\Sigma}^1 & \dots & \hat{\Sigma}^{1E} \\ \vdots & \ddots & \vdots \\ \hat{\Sigma}^{E1} & \dots & \hat{\Sigma}^E \end{bmatrix} \right) \quad (27)$$

$$882 =: \hat{S} \sim \mathcal{N}(\hat{\mu}, \hat{\Sigma})$$

883 with  $\hat{\Sigma}^{ef} := \text{Cov}[\hat{S}_{t+1}^e, \hat{S}_{t+1}^f]$ . We aim to track  $S_{t+1}$  such that

$$884 HS_{t+1} \sim \mathcal{N}(\hat{\mu}, \hat{\Sigma}) \quad (28)$$

885 where we use  $H = [I, I, \dots, I]^\top \in \mathbb{R}^{n_S \cdot E \times n_S}$  to project  $S_{t+1}$  to the dimension of the joint  $\hat{S}$ .

886 We define the maximum likelihood loss

$$887 \mathcal{L}(S_{t+1}) = p(\hat{S} | S_{t+1}) = \frac{1}{|2\pi\hat{\Sigma}|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\hat{S} - HS_{t+1}) \hat{\Sigma}^{-1} (\hat{S} - HS_{t+1}) \right) \quad (29)$$

888 such that

$$889 \log(\mathcal{L}(S_{t+1})) = -\frac{1}{2} \log(|2\pi\hat{\Sigma}|) - \frac{1}{2} (\hat{S} - HS_{t+1}) \hat{\Sigma}^{-1} (\hat{S} - HS_{t+1}). \quad (30)$$

890 We aim to obtain the maximizer of the log-likelihood such that

$$891 \bar{S}_{t+1} = \arg \max_{S_{t+1}} \log(\mathcal{L}(S_{t+1})). \quad (31)$$

892 Consequently,

$$893 \frac{\partial}{\partial S_{t+1}} \log(\mathcal{L}(S_{t+1})) = -\frac{1}{2} H^\top \hat{\Sigma}^{-1} (\hat{S} - HS_{t+1}) := 0$$

$$894 \Rightarrow H^\top \hat{\Sigma}^{-1} \hat{S} - H^\top \hat{\Sigma}^{-1} H \bar{S}_{t+1} = 0 \quad (32)$$

$$895 \Rightarrow \bar{S}_{t+1} = \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{S}.$$

896 As a result, we obtain

$$897 \bar{\mu} = \mathbb{E}[\bar{S}_{t+1}] = \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{\mu} \quad (33)$$

918 and

$$\begin{aligned}
919 \quad \bar{\Sigma} &= \text{Var} [\bar{S}_{t+1}] = \text{Var} \left[ \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{S} \right] \\
920 &= \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \text{Var} [\hat{S}] \hat{\Sigma}^{-1} H \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} \\
921 &= \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{\Sigma} \hat{\Sigma}^{-1} H \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} = \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1}
\end{aligned} \tag{34}$$

926 which corresponds to standard results in Kalman fusion. However, as the cross-correlations  $\hat{\Sigma}^{ef}$  are unknown in practice, we approximate the Kalman fusion results (33) and (34) using covariance intersection fusion Julier & Uhlmann (2001) with uniform weights, making use of Assumption 1. This results in

$$\bar{\Sigma} = \left( \frac{1}{E} \sum_{e=1}^E \left( \hat{\Sigma}^e \right)^{-1} \right)^{-1} \tag{35}$$

933 and

$$\bar{\mu} = \bar{\Sigma} \left( \frac{1}{E} \sum_{e=1}^E \left( \hat{\Sigma}^e \right)^{-1} \hat{\mu}^e \right). \tag{36}$$

937 Hence, we can estimate the environment state as

$$\bar{S}_{t+1} = \bar{\mu}(S_t, A_t) + \bar{L}(S_t, A_t)W_t, \tag{37}$$

940 with  $\bar{L}\bar{L}^\top = \bar{\Sigma}$  and  $W_t \sim \mathcal{N}(0, I)$ .  $\square$

#### 942 D.2.2 PROOF OF LEMMA 2

943 *Proof.* We continue here using the quantities we estimated in the previous section. To estimate  $\Sigma^\Delta$ , we interpret  $\{\hat{\mu}^1\}_{e=1}^E$  as samples from a distribution whose mean is known to be  $\bar{\mu}$ . With this, the maximum likelihood estimate of  $\Sigma^\Delta$  can be obtained trivially as

$$\bar{\Sigma}^\Delta = \frac{1}{E} \sum_{e=1}^E (\hat{\mu}^e - \bar{\mu})(\hat{\mu}^e - \bar{\mu})^\top. \tag{38}$$

950  $\square$

#### 952 D.3 INFOPROP STATE

953 As introduced in (15), the Infoprop state is defined as

$$\tilde{S}_{t+1} := \mathbb{E} \left[ \bar{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1} \right] = \mathbb{E}_{\tilde{P}_{S,IP}} \left[ S_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}, U_t \right] \tag{39}$$

957 Combining (11) and Assumption 2, we have

$$\hat{S}_{t+1} = \check{S}_{t+1} + L^\Delta(S_t, A_t)N_t. \tag{40}$$

960 Plugging the respective maximum likelihood estimates into (40) yields

$$\hat{S}_{t+1} = \bar{S}_{t+1} + \bar{L}^\Delta(S_t, A_t)N_t \tag{41}$$

963 with

$$\bar{S}_{t+1} = \bar{\mu}(S_t, A_t) + \bar{L}(S_t, A_t)W_t \tag{42}$$

965 according to (13). As we can generally consider model uncertainty as independent from process noise, i.e.  $N_t \perp W_t$ , the Infoprop state  $\tilde{S}_{t+1} = \mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}]$  can be computed using a standard Kalman update.

968 The general form of the Kalman update Simon (2006) considers two Gaussian RVs  $X \sim \mathcal{N}(\mu_X, \Sigma_X)$  and  $Y = X + N$  with  $N \sim \mathcal{N}(0, \Sigma_N)$  and  $X \perp N$ . Then, given an observation  $y$  we can compute the conditional expectation of  $X$

$$\mathbb{E}[X|Y=y] \sim \mathcal{N}(\mu_{X|Y=y}, \Sigma_{X|Y=y}) \tag{43}$$

972 with

$$973 \mu_{X|Y=y} = \mu_X + K(y - \mu_X), \quad (44)$$

$$974 \Sigma_{X|Y=y} = (I - K) \Sigma_X, \quad (45)$$

975 and

$$976 K = \Sigma_X (\Sigma_X + \Sigma_N)^{-1}. \quad (46)$$

977 Following (15), we can compute the Infoprop state via (43) choosing

$$978 \mu_X = \bar{\mu}(s_t, a_t), \quad (47)$$

$$980 \Sigma_X = \bar{\Sigma}(s_t, a_t), \quad (48)$$

$$981 \Sigma_N = \bar{\Sigma}^\Delta(s_t, a_t), \quad (49)$$

982 and

$$983 y = \bar{\mu}(s_t, a_t) + \bar{L}(s_t, a_t)w_t + \bar{L}^\Delta(s_t, a_t)n_t. \quad (50)$$

984 This yields the propagation equation of the Infoprop state

$$985 \tilde{S}_{t+1} = \tilde{\mu}(S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}) + \tilde{L}(S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1})U_t \quad (51)$$

986 with

$$987 \tilde{\mu}(s_t, a_t, \hat{s}_{t+1}) = \bar{\mu}(s_t, a_t) + K(s_t, a_t) (\bar{L}(s_t, a_t)w_t + \bar{L}^\Delta(s_t, a_t)n_t), \quad (52)$$

$$988 \tilde{\Sigma}(s_t, a_t, \hat{s}_{t+1}) = (I - K(s_t, a_t)) \bar{\Sigma}(s_t, a_t), \quad (53)$$

$$989 K(s_t, a_t) = \bar{\Sigma}(s_t, a_t) (\bar{\Sigma}(s_t, a_t) + \bar{\Sigma}^\Delta(s_t, a_t))^{-1}, \quad (54)$$

$$990 \tilde{L}(s_t, a_t) \tilde{L}(s_t, a_t)^\top = \tilde{\Sigma}(s_t, a_t), \quad (55)$$

991 and

$$992 \bar{L}^\Delta(s_t, a_t) \bar{L}^\Delta(s_t, a_t)^\top = \bar{\Sigma}^\Delta(s_t, a_t). \quad (56)$$

#### 993 D.4 INDUCED STATE DISTRIBUTION BY THE INFOPROP ROLLOUT

994 **Lemma 3.** *As introduced in (57), the next state distribution induced by the Infoprop rollout is the same as that given by the estimated ground truth:*

$$995 \tilde{S}_{t+1} \stackrel{\text{dist}}{=} \bar{S}_{t+1} \quad (57)$$

996 *Proof.* We show equality in distribution via comparison of the cumulative distribution functions (CDF) of  $\tilde{S}_{t+1}$  and  $\bar{S}_{t+1}$ . If we can show that the CDFs are identical, i.e.  $\mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) \quad \forall \bar{s}_{t+1} \in \mathcal{S}$ , the equality in distribution follows.

1000 We compute  $\mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1})$  using  $\tilde{S}_{t+1} = \mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}]$  and marginalizing over  $\hat{S}_{t+1}$

$$1001 \mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \int_{\mathcal{S}} \mathbb{P}(\mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1}] \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) f_{\hat{S}_{t+1}}(\hat{s}_{t+1}) d\hat{s}_{t+1} \quad (58)$$

1002 with  $f_{\hat{S}_{t+1}}$  the probability density function of  $\hat{S}_{t+1}$ .

1003 By construction,  $\mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1}]$  describes the behavior of  $\bar{S}_{t+1}$  given  $\hat{S}_{t+1}$ . Consequently,

$$1004 \mathbb{P}(\mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1}] \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) = \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) \quad (59)$$

1005 and therefore

$$1006 \mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \int_{\mathcal{S}} \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) f_{\hat{S}_{t+1}}(\hat{s}_{t+1}) d\hat{s}_{t+1}. \quad (60)$$

1007 The right hand side of (60) represents the law of total probability for  $P(\bar{S}_{t+1} \leq \bar{s}_{t+1})$

$$1008 \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) = \int_{\mathcal{S}} \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) f_{\hat{S}_{t+1}}(\hat{s}_{t+1}) d\hat{s}_{t+1}. \quad (61)$$

1009 Therefore, we have

$$1010 \mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) \quad \forall \bar{s}_{t+1} \in \mathcal{S} \quad (62)$$

1011 and can conclude

$$1012 \tilde{S}_{t+1} \stackrel{\text{dist}}{=} \bar{S}_{t+1}. \quad (63)$$

1013  $\square$

## D.5 INFORMATION LOSS ALONG A INFOPROP ROLLOUT

**Lemma 4.** *As introduced in (17), the total information loss incurred during a Infoprop equals the accumulated entropy of the Infoprop state:*

$$\mathbb{H}\left(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T\right) = \sum_{t=0}^{T-1} \mathbb{H}\left(\tilde{S}_{t+1}\right) \quad (64)$$

*Proof.*

$$\begin{aligned} & \mathbb{H}\left(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T\right) \\ &= \sum_{t=0}^{T-1} \mathbb{H}\left(\bar{S}_{t+1} | \bar{S}_1, \bar{S}_2, \dots, \bar{S}_t, S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T\right) \\ &\stackrel{(a)}{=} \sum_{t=0}^{T-1} \mathbb{H}\left(\bar{S}_{t+1} | \bar{S}_1, \bar{S}_2 \dots \bar{S}_t, S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_t\right) \\ &\stackrel{(b)}{=} \sum_{t=0}^{T-1} \mathbb{H}\left(\bar{S}_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}\right) \\ &= \sum_{t=0}^{T-1} \mathbb{H}\left(\tilde{S}_{t+1}\right) \end{aligned} \quad (65)$$

where (a) follows from causality and (b) follows from the Markov property.  $\square$

## E EXPERIMENTS

### E.1 EXPERIMENTAL SETUP

We used Weights&Biases<sup>4</sup> for logging our experiments and run 5 random seeds per experiment.

The respective hyperparameters for Infoprop-Dyna on MuJoCo are given below. Table 2 addresses model learning, Table 3 the Infoprop mechanism, and Table 4 training the model-free agent.

Table 2: Hyperparameters used to train the model of Infoprop-Dyna in the Mujoco Tasks.

| Hyperparameter              | Halfcheetah           | Walker | Hopper | Ant     |
|-----------------------------|-----------------------|--------|--------|---------|
| ensemble size $E$           | 7                     |        |        |         |
| number of hidden neurons    | 200                   |        |        | 400     |
| number of hidden layers     | 4                     |        |        |         |
| learning rate               | 0.0003                | 0.0006 | 0.0004 | 0.001   |
| weight decay                | 0.00005               | 0.0007 | 0.0008 | 0.00002 |
| patience for early-stopping | 10                    | 9      | 8      | 9       |
| retrain interval            | 250 environment steps |        |        |         |

<sup>4</sup><https://wandb.ai/site>

Table 3: Hyperparameters of the Infoprop rollouts in the Mujoco Tasks.

| Hyperparameter                            | Halfcheetah           | Walker | Hopper | Ant |
|---|-----------------------|--------|--------|-----|
| accurate quantile $\zeta_1$               | 0.99                  |        |        |     |
| exceptionally accurate quantile $\zeta_2$ | 0.01                  |        |        |     |
| scaling factor $\xi$                      | 100                   |        |        |     |
| rollout interval                          | 250 environment steps |        |        |     |
| rollout batch size                        | 100000                |        |        |     |

Table 4: Hyperparameters used to train the SAC agent of Infoprop-Dyna in the Mujoco Tasks.

| Hyperparameter           | Halfcheetah | Walker | Hopper | Ant    |
|--------------------------|-------------|--------|--------|--------|
| number of hidden neurons | 1024        |        | 512    | 1024   |
| number of hidden layers  | 2           |        |        |        |
| learning rate            | 0.0003      | 0.0002 | 0.0004 | 0.0005 |
| SAC target entropy       | -6          | -7     | 1      | 0      |
| target update interval   | 1           | 4      | 6      | 5      |
| update steps $G$         | 10          |        |        | 20     |

The results for SAC, MBPO and MACURA are obtained from Frauenknecht et al. (2024).

## E.2 PREDICTION QUALITY

We provide additional results for the rollout consistency experiments introduced in Section 6.2. Figure 6 depicts model-based rollouts for the 10<sup>th</sup> dimension of hopper under MBPO, MACURA and Infoprop-Dyna when setting the maximum rollout length of all approaches to 100. In the original experiment depicted in Figure 4b the maximum rollout length was 11 for MBPO and 10 for MACURA, following the hyperparameter settings reported in the respective publications Janner et al. (2019); Frauenknecht et al. (2024).

We observe a vastly spread distribution of MBPO rollouts, as every rollout is propagated for 100 steps, irrespective of model uncertainty, as long as it does not reach a terminal state of the hopper task. MACURA rollouts have an improved consistency compared to MBPO, especially in the beginning of the rollouts. Over long horizons, however, the TS propagation mechanism and the single-step termination criterion cannot produce consistent data. In contrast, Infoprop-Dyna is able to propagate consistent rollouts over long horizons.

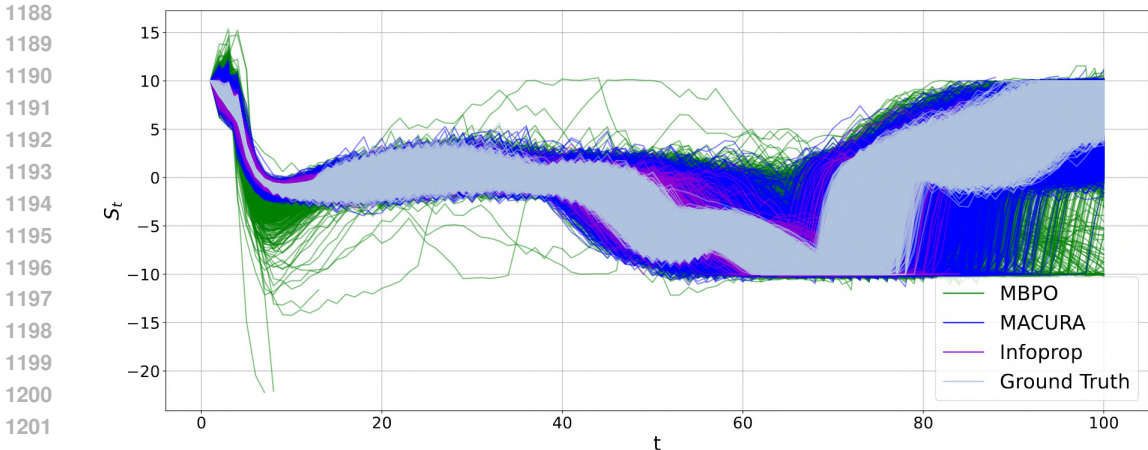


Figure 6: Rollout consistency MBPO vs. MACURA vs. Infoprop-Dyna for 100 steps. Comparison of the respective rollout mechanisms similar to Figure 4b but with a maximum rollout length of 100 for all approaches.

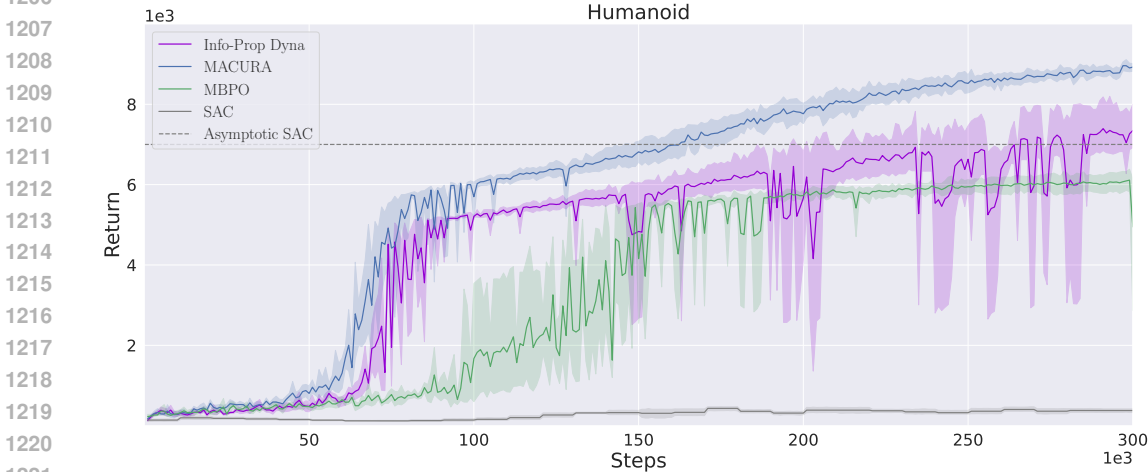


Figure 7: Performance on Humanoid

### E.3 PERFORMANCE ON HUMANOID

Figure 7 depicts the return on MuJoCo humanoid. We observe instabilities in the performance of Infoprop-Dyna towards the end of training. We assume this occurs due to overfitting and plasticity loss in the critic of the model-free learner Nikishin et al. (2022); D’Oro et al. (2023). This is reflected in the peaking critic loss depicted in Figure 8 concurrently with the performance drops. We set the update ratio  $G$  (see Algorithm 3) to a relatively low value of 10 which explains the slower learning behavior than MACURA. For higher values of  $G$ , instabilities occur even earlier in the training process, underscoring our assumption of overfitting critics.

Model rollout inconsistency does not appear to be the destabilizing factor, as rollout data is consistent with the environment distribution as depicted in Figure 10 and the rollout adaption mechanism seems to react to policy shifts induced by high critic losses through reducing the average rollout length as depicted in Figure 10.

### E.4 INVESTIGATING INSTABILITIES IN LEARNING

Although Infoprop gives better quality data over longer rollout horizons than TS rollouts, we observe instabilities in learning when naively integrating Infoprop into the conventional Dyna setting. We hypothesize that the main cause of these instabilities is due to the agent overfitting to the higher

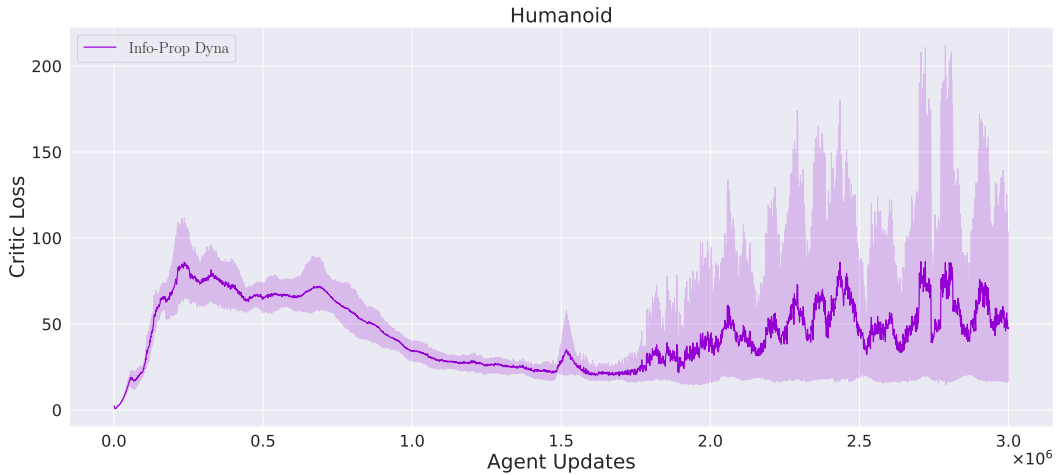
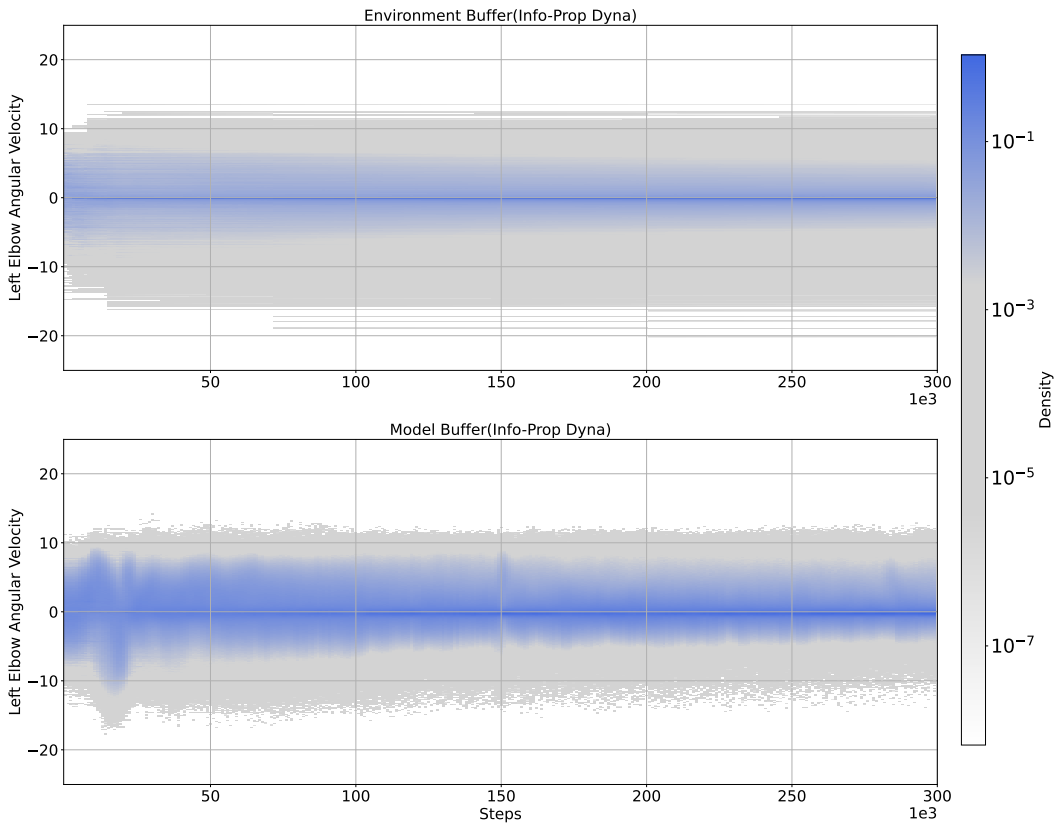


Figure 8: Critic Loss on Humanoid

Figure 9: Comparison between  $\mathcal{D}_{env}$  and  $\mathcal{D}_{mod}$  for the 45<sup>th</sup> dimension of Humanoid

1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293

quality data produced by Infoprop rollouts, followed by loss of plasticity Nikishin et al. (2022); D’Oro et al. (2023). To investigate this, we carried out an ablation by varying the values of  $\zeta_1$ , which we introduced in Equation 22. This hyperparameter controls the size of the subset  $\mathcal{E}$  where the model is considered sufficiently accurate. The smaller the value of  $\zeta_1$ , the more aggressive the filtering of single-step information losses, leading to a smaller  $\mathcal{E}$ .

1294  
1295

Figure 11 shows the returns obtained on the Hopper task for three values of  $\zeta_1$ . For  $\zeta_1 = 0.97$ , we see that the returns are unstable throughout training, even though this setting gives the best quality data. On the other hand,  $\zeta_1 = 0.9999$  produces a more stable learning curve compared to



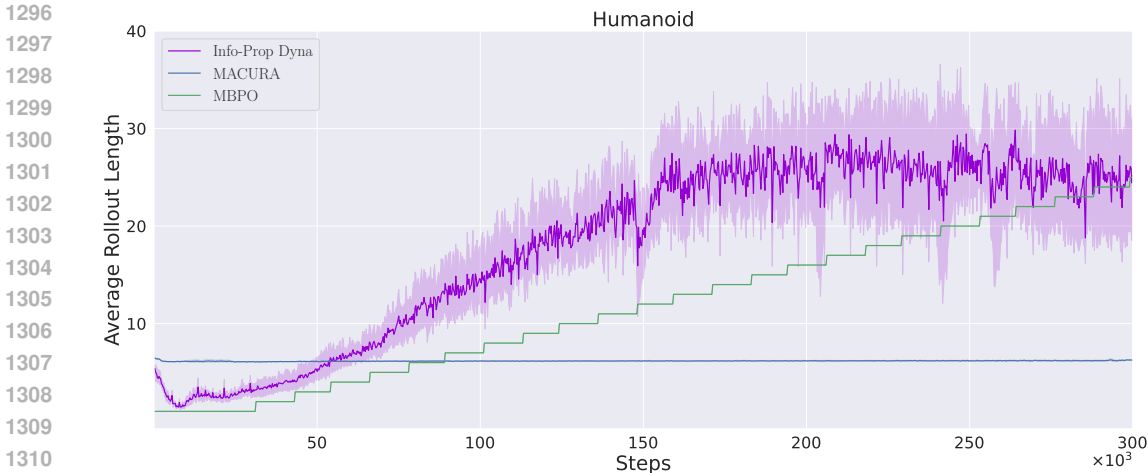


Figure 10: Average Rollout Length on Humanoid

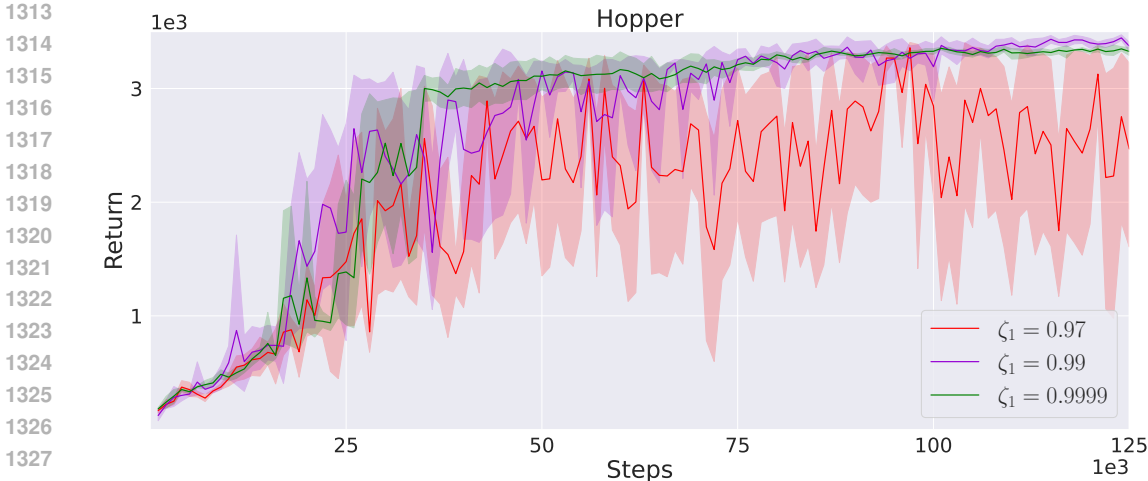


Figure 11: Ablation Study on Hopper.

1331  $\zeta_1 = 0.99$ , which was used for all the experiments in Section 6. This shows that better data quality  
 1332 does not necessarily lead to better training performance since if that was the case,  $\zeta_1 = 0.97$  would  
 1333 have produced the best performance. A similar observation is reported in Frauenknecht et al. (2024),  
 1334 where low values of the scaling factor  $\xi$ , corresponding to accurate model rollouts, led to instabilities  
 1335 in learning.

1336 Our observations show that producing high-quality synthetic data in the conventional Dyna setting  
 1337 leads to issues seen in MFRL when using a high update-to-data (UTD) ratio. There have been  
 1338 recent works on regularization methods to counteract agent overfitting and loss of plasticity. One  
 1339 such approach is applying layer normalization Smith et al. (2023); Nauman et al. (2024). Figure 12  
 1340 shows the same settings as in Figure 11 but with layer normalization applied to the critic and actor  
 1341 networks. It can be seen that even for  $\zeta_1 = 0.97$ , the learning is stable.

1342 The primary aim of this paper is to introduce the conceptual framework of the Infoprop rollout,  
 1343 as well as show its application to MBRL. Hence, we do not spend additional effort on tuning the  
 1344 hyperparameters or adding regularizations since this takes us away from the main objective. We  
 1345 defer such enhancements for future work.

1346  
 1347  
 1348  
 1349

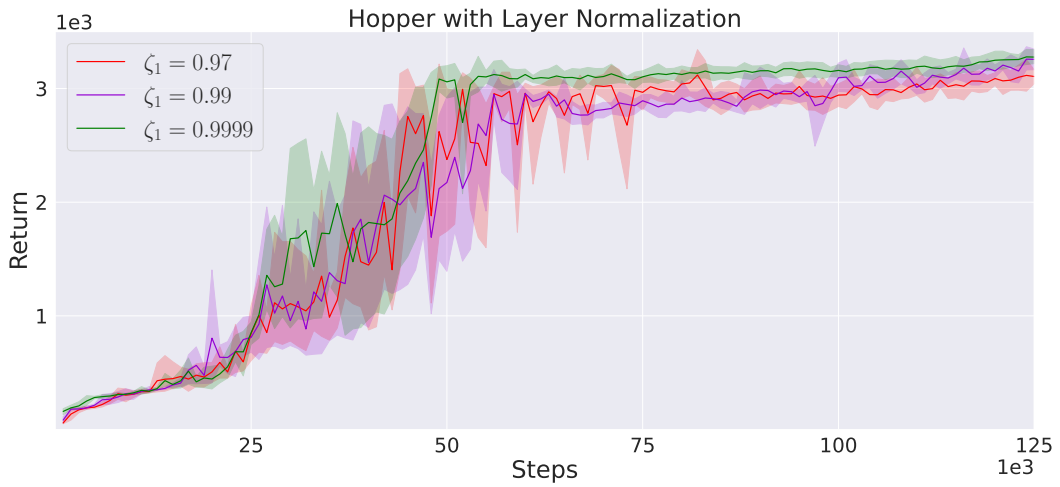


Figure 12: Ablation Study on Hopper with layer normalization.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

# ON ROLLOUTS IN MODEL-BASED REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model-based reinforcement learning (MBRL) seeks to enhance data efficiency by learning a model of the environment and generating synthetic rollouts from it. However, accumulated model errors during these rollouts can distort the data distribution, negatively impacting policy learning and hindering long-term planning. Thus, the accumulation of model errors is a key bottleneck in current MBRL methods. We propose *Infoprop*, a model-based rollout mechanism that separates aleatoric from epistemic model uncertainty and reduces the influence of the latter on the data distribution. Further, *Infoprop* keeps track of accumulated model errors along a model rollout and provides termination criteria to limit data corruption. We demonstrate the capabilities of *Infoprop* in the *Infoprop-Dyna* algorithm, reporting state-of-the-art performance in Dyna-style MBRL on common MuJoCo benchmark tasks while substantially increasing rollout length and data quality.

## 1 INTRODUCTION

Reinforcement learning (RL) has emerged as a powerful framework for solving complex decision-making tasks like racing Vasco et al. (2024); Kaufmann et al. (2023) and gameplay OpenAI et al. (2019); Bi & D’Andrea (2024). However, when applying RL in real-world scenarios, a significant challenge is data inefficiency, which hinders the practicality of standard RL methods. Model-based reinforcement learning (MBRL) addresses this issue by learning an internal model of the environment Deisenroth & Rasmussen (2011); Chua et al. (2018); Janner et al. (2019); Hafner et al. (2020). By generating simulated interactions through model rollouts, MBRL can make informed decisions while substantially reducing the need for real-world data collection.

The quality of data from model-based rollouts is critical for MBRL performance. Model errors can distort the data distribution and hurt policy learning. Long-horizon planning is desirable, however, often infeasible as model errors accumulate over time. This effect is demonstrated in Figure 1. Even for a simple toy example (described in Appendix B), we see the data distribution of model-based rollouts under the state-of-the-art Trajectory Sampling (TS) Chua et al. (2018) scheme diverging quickly from the ground truth distribution of environment rollouts. Thus, data from TS rollouts can even be harmful to policy learning after a couple of time steps. This is largely because the TS mechanism does not explicitly address the effect of model errors on the propagated data distribution.

To tackle this challenge, we propose *Infoprop*, a novel model-based rollout mechanism that mitigates data distortion by addressing two key questions: *How to propagate?* and *When to stop?* We build our mechanism on explicitly leveraging the ability of common MBRL models to distinguish between aleatoric uncertainty due to process noise and epistemic uncertainty due to lack of data Lakshminarayanan et al. (2017); Becker & Neumann (2022). Making use of this property leads to substantially improved data consistency as depicted in Figure 1. In particular, we

- estimate and remove the stochasticity due to model error from the predictive distribution;
- formulate stopping criteria based on information loss to limit error accumulation; and
- demonstrate the potential of *Infoprop* as a direct plugin to standard MBRL methods using the example of Dyna-style MBRL. The resulting *Infoprop-Dyna* algorithm yields state-of-the-art performance in MBRL on common MuJoCo tasks, while substantially improving the data consistency of model-based rollouts and thus allowing for longer rollout horizons.

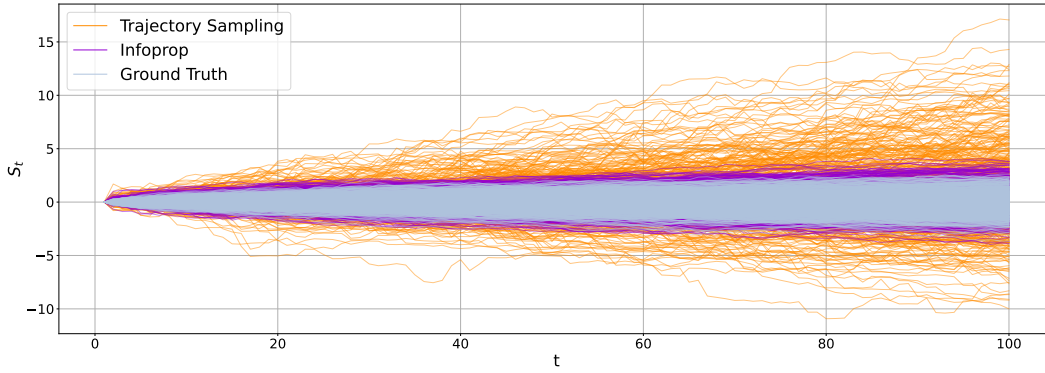


Figure 1: Comparing Data Consistency of Model-based Rollouts. *Trajectories under the proposed Infoprop mechanism follow the ground truth distribution of environment rollouts closely while rolling out the same model under the common TS scheme Chua et al. (2018) results in distorted data.*

## 2 BACKGROUND

In the following, we introduce the fundamental concepts of information theory and MBRL. [Appendix A provides an overview of the notation introduced and used in the remainder of the paper.](#)

### 2.1 INFORMATION THEORY

We will estimate the degree of data corruption in Infoprop rollouts using information-theoretic arguments. Information theory serves to quantify the uncertainty of a random variable (RV) Shannon (1948). Given the discrete RVs  $X : \Omega \rightarrow \mathcal{X}$  and  $Y : \Omega \rightarrow \mathcal{Y}$ , the marginal entropy  $\mathbb{H}(X) = -\sum_{x \in \mathcal{X}} \mathbb{P}[X = x] \log_2(\mathbb{P}[X = x])$  describes the average uncertainty about  $X$  in bits. Further, the conditional entropy  $\mathbb{H}(X|Y = y) = -\sum_{x \in \mathcal{X}} \mathbb{P}[X = x|Y = y] \log_2(\mathbb{P}[X = x|Y = y])$  gives the uncertainty about  $X$ , given a realization of  $Y$ . Based on marginal and conditional entropy, the reduction in uncertainty about  $X$  given a realization of  $Y$  is described by mutual information

$$\mathbb{I}(X; Y = y) = \mathbb{H}(X) - \mathbb{H}(X|Y = y), \quad (1)$$

with  $\mathbb{I}(X; Y = y) = 0$  if the RVs are independent. In the following, we focus on Gaussian RVs and use the notion of quantized entropy Cover & Thomas (2006). ~~For a RV  $Z \in \mathcal{Z} \subseteq \mathbb{R}^{n_Z}$  with  $Z \sim \mathcal{N}(\mu_Z, \Sigma_Z)$  and discretization step size  $\Delta z^{(k)}$  of the  $k^{\text{th}}$  dimension, the quantized entropy is~~

$$\mathbb{H}(Z) = \frac{1}{2} \log_2((2\pi e)^{n_Z} |\Sigma_Z|) - \sum_{k=1}^{n_Z} \log_2(\Delta z^{(k)}).$$

[with details provided in Appendix D.1.](#)

### 2.2 REINFORCEMENT LEARNING

Reinforcement learning addresses sequential decision-making problems where the environment is typically modeled as a discrete-time Markov decision process (MDP) represented by the tuple  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{R}, P_{\mathcal{R}}, P_{\mathcal{S}}, \xi_0, \gamma\}$ . Here,  $\mathcal{S} \subseteq \mathbb{R}^{n_S}$  denotes the state space with  $S_t \in \mathcal{S}$  being the RV of the state at time  $t$  and  $s_t$  its realization. Similarly,  $\mathcal{A} \subseteq \mathbb{R}^{n_A}$  represents the action space with  $A_t \in \mathcal{A}$  the RV and  $a_t$  the realization of the action as well as  $\mathcal{R} \subseteq \mathbb{R}$  the set of rewards with  $R_t \in \mathcal{R}$  and  $r_t$  the reward at time  $t$ . We make the common simplifying assumption Bellemare et al. (2023) that the next state and reward are independent given the current state-action pair. Thus, a transition step in the environment can be expressed concerning a reward kernel  $P_{\mathcal{R}}$  and a dynamics kernel  $P_{\mathcal{S}}$  as

$$R_{t+1} \sim P_{\mathcal{R}}(\cdot|S_t, A_t) \quad \text{and} \quad S_{t+1} \sim P_{\mathcal{S}}(\cdot|S_t, A_t). \quad (2)$$

Further, initial states are distributed according to  $S_0 \sim \xi_0$ , and actions according to the policy  $A_t \sim \pi(\cdot|S_t)$ . We aim to learn an optimal policy  $\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t R_{t+1}]$  that maximizes the expected sum of rewards discounted by  $\gamma \in [0, 1)$ , referred to as return.

## 2.3 MODEL-BASED REINFORCEMENT LEARNING

There are four main categories of MBRL that all build on model-based rollouts. (i) Dyna-style methods Sutton (1991); Janner et al. (2019) use model-based rollouts to generate training data for a model-free agent. (ii) Model-based planning approaches Chua et al. (2018); Williams et al. (2017); Nagabandi et al. (2018); Hafner et al. (2019) do not learn an explicit policy but perform planning via model rollouts during deployment. (iii) Analytic gradient methods Deisenroth & Rasmussen (2011); Hafner et al. (2020; 2021; 2023) optimize the policy by backpropagating the performance gradient through model rollouts. (iv) Value-expansion approaches Feinberg et al. (2018); Buckman et al. (2018) stabilize the temporal difference target using model-based rollouts.

The model architecture of an MBRL algorithm determines the set of mechanisms for model rollouts. In this work, we focus on rolling out the particularly successful class of aleatoric epistemic separator (AES) models Lakshminarayanan et al. (2017); Becker & Neumann (2022) that can distinguish aleatoric uncertainty corresponding to the estimate of process noise from epistemic uncertainty.

## 2.4 ENVIRONMENT INTERACTION VS. MODEL-BASED ROLLOUTS

Model-based rollouts aim to substitute environment interaction in MBRL. Thus, we compare the data generation process through environment interaction to the process of model-based rollouts.

We model environment dynamics as a nonlinear function  $\mu(S_t, A_t)$  with additive heteroscedastic process noise that is normally distributed with variance  $\Sigma(S_t, A_t)$ . Thus, environment rollouts, as depicted in Figure 1, are generated by iterating the dynamics

$$S_{t+1} = \mu(S_t, A_t) + L(S_t, A_t)W_t, \quad (3)$$

with  $L(S_t, A_t)L(S_t, A_t)^\top = \Sigma(S_t, A_t)$  and the process noise  $W_t \sim \mathcal{N}(0, I)$ . Consequently, the transition kernel<sup>1</sup> of the environment is defined as  $P_S(\cdot | S_t, A_t) = \mathcal{N}(\mu(S_t, A_t), \Sigma(S_t, A_t))$ .

In MBRL, however, we do not have access to  $P_S$  directly but typically rely on a parametric model with the random parameters  $\Theta_t \in \vartheta$ . Besides estimates of nonlinear dynamics  $\hat{\mu}_{\Theta_t}(S_t, A_t)$  and process noise  $\hat{\Sigma}_{\Theta_t}(S_t, A_t)$ , AES models provide an estimate of the parameter distribution  $\Theta_t \sim \mathbb{P}_\Theta$ , e.g. via ensembling Lakshminarayanan et al. (2017) or dropout Becker & Neumann (2022). These models are typically propagated using the TS Chua et al. (2018) rollout mechanism via iterating

$$S_{t+1} = \hat{\mu}_{\Theta_t}(S_t, A_t) + \hat{L}_{\Theta_t}(S_t, A_t)W_t \quad (4)$$

with  $\hat{L}_{\Theta_t}(S_t, A_t)\hat{L}_{\Theta_t}(S_t, A_t)^\top = \hat{\Sigma}_{\Theta_t}(S_t, A_t)$ ,  $W_t \sim \mathcal{N}(0, I)$ , and  $\Theta_t \sim \mathbb{P}_\Theta$ . This results in the TS rollouts in Figure 1 and induces the kernel  $\hat{P}_{S,TS}(\cdot | S_t, A_t) = \mathcal{N}(\hat{\mu}_{\Theta_t}(S_t, A_t), \hat{\Sigma}_{\Theta_t}(S_t, A_t))$ .

The majority of recent MBRL approaches use the TS rollout mechanism, e.g.

[Chua et al. \(2018\)](#); [Buckman et al. \(2018\)](#); [Janner et al. \(2019\)](#); [Pan et al. \(2020\)](#); [Yu et al. \(2020\)](#); [Luis et al. \(2023\)](#)  
[Chua et al. \(2018\)](#); [Becker & Neumann \(2022\)](#); [Janner et al. \(2019\)](#); [Pan et al. \(2020\)](#); [Yu et al. \(2020\)](#); [Luis et al. \(2023\)](#)  
 . Pseudocode is provided in Algorithm 2 of Appendix C.

## 3 PROBLEM STATEMENT

Revisiting Figure 1 allows us to illustrate the effects of different sources of stochasticity by comparing environment interaction under  $P_S$  to TS rollouts under  $\hat{P}_{S,TS}$ . While different realizations of process noise  $w_t \sim \mathcal{N}(0, I)$  allow for keeping track of the environment distribution, the sampling process  $\theta_t \sim \mathbb{P}_\Theta$  introduces additional stochasticity that leads to an overestimated total variance in the TS rollout distribution. This effect is amplified through the continued propagation of erroneous predictions making data at later steps unfit for policy learning. We ask the following questions:

- (i) How can we construct a predictive distribution closely resembling environment dynamics?
- (ii) How can we quantify the degree of data corruption due to model error?
- (iii) When should model-based rollouts be terminated due to data corruption?

<sup>1</sup>As  $P_R$  typically is a known deterministic function in the context of MBRL, while  $P_S$  is the unknown object we aim to model, the discussion henceforth focuses on approximating  $P_S$  without loss of generality.

We address these questions by proposing the *Infoprop* rollout mechanism. We combine different beliefs under  $\mathbb{P}_\Theta$  Infoprop isolates and removes epistemic uncertainty for an improved predictive distribution, keep-keeps track of data corruption using information-theoretic arguments, and terminate terminates rollouts based on the degree of corruption.

## 4 INFOPROP ROLLOUT MECHANISM

In the following, we introduce the *Infoprop* mechanism for model-based rollouts. As depicted in Figure 2, we decompose model predictions into a signal fraction representing the environment dynamics and noise fraction introduced by model error. This perspective allows to interpret model rollouts as communication through a noisy channel. We estimate both the signal and the noise distribution and use these to infer a belief over the environment state, given an observation of the model state. This belief state represents the foundation of the Infoprop rollout mechanism.

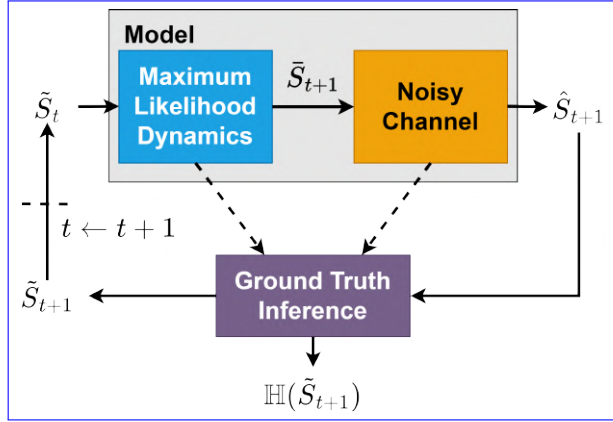


Figure 2: Infoprop block diagram

### 4.1 THEORETICAL SETUP

First, we introduce additional notation to specify RVs under different transition kernels. We define

$$\check{S}_{t+1} := \mathbb{E}_{P_S} [S_{t+1} | S_t = s_t, A_t = a_t, W_t]$$

**Definition 1** (Environment state). We define the environment state as the conditional expectation under environment dynamics over the next state, the environment dynamics given a realization of a state-action pair :-

$$\check{S}_{t+1} := \mathbb{E}_{P_S} [S_{t+1} | S_t = s_t, A_t = a_t, W_t]. \quad (5)$$

Thus,  $\check{S}_{t+1}$  is an RV, where the randomness is induced by the process noise and has an aleatoric nature. If we additionally condition on the realization  $W_t = w_t$ , we obtain a deterministic object. Similarly, we introduce

**Definition 2** (Model state). We define the model state as the conditional expectation under the TS kernel as  $\hat{P}_{S,TS}$

$$\hat{S}_{t+1} := \mathbb{E}_{\hat{P}_{S,TS}} [S_{t+1} | S_t = s_t, A_t = a_t, W_t, \Theta_t]. \quad (6)$$

As discussed in Section 3, stochasticity in  $\hat{S}_{t+1}$  is induced not only by  $W_t$  but also by the randomness in the parameters  $\Theta_t$ . We project the uncertainty in the parameter space  $\vartheta$  to  $\mathcal{S}$  via an error process.

**Definition 3** (Model error process). We define a model error process  $\Delta_t$

$$\Delta_t = \hat{S}_{t+1} - \check{S}_{t+1} \quad (7)$$

that, given a realization of process noise  $W_t = w_t$ , projects stochasticity due to  $\Theta_t$  into uncertainty in  $\vartheta$  to  $\mathcal{S}$ , such that,

$$\mathbb{E} [\Delta_t | W_t = w_t] = \mathbb{E}_{\hat{P}_{S,TS}} [S_{t+1} | s_t, a_t, w_t, \Theta_t] - \mathbb{E}_{P_S} [S_{t+1} | s_t, a_t, w_t]. \quad (8)$$

which we refer to as epistemic uncertainty. The parameter distribution  $\Theta_t \sim \mathbb{P}_\Theta$  can induce arbitrarily complex distributions  $\Delta_t \sim \mathbb{P}_\Delta$ . To simplify the analysis, we solely consider the first two moments of  $\mathbb{P}_\Delta$ , such that

$$\hat{S}_{t+1} = \check{S}_{t+1} + \mu^\Delta(S_t, A_t) + L^\Delta(S_t, A_t) N_t$$

with  $L^\Delta(S_t, A_t)L^\Delta(S_t, A_t)^\top = \Sigma^\Delta(S_t, A_t)$  and  $N_t \sim \mathcal{N}(0, I)$ . We refer to  $\mu^\Delta(S_t, A_t)$  as the model bias,  $\Sigma^\Delta(S_t, A_t)$  the epistemic variance, and  $N_t$  the epistemic noise.

the projected parameter uncertainty as epistemic uncertainty.

Further, we restrict model usage to a sufficiently accurate subset  $\mathcal{E} \subseteq \mathcal{S} \times \mathcal{A}$ , as proposed in Frauenknecht et al. (2024). We define  $\mathcal{E}$  amenable to the Infoprop setting in Section 4.4 and make the following assumptions when performing model-based rollouts in  $\mathcal{E}$ :

**Assumption 1** (Consistent estimator of aleatoric uncertainty). *The model’s predictive variance  $\hat{\Sigma}_{\Theta_t}$  is a consistent estimator of  $\Sigma$  following the definition of Julier & Uhlmann (2001), i.e.  $(\hat{\Sigma}_{\Theta_t}(S_t, A_t) - \Sigma(S_t, A_t)) \succcurlyeq 0 \quad \forall (S_t, A_t) \in \mathcal{E}$ .*

$$\left( \hat{\Sigma}_{\Theta_t}(S_t, A_t) - \Sigma(S_t, A_t) \right) \succcurlyeq 0 \quad \forall (S_t, A_t) \in \mathcal{E}. \quad (9)$$

**Assumption 2** (Unbiased estimator). *The model bias  $\mu^\Delta$  is negligible. Thus  $\hat{S}_{t+1}$  according to (4) is an unbiased estimator of  $\check{S}_{t+1}$  according to (3), i.e.  $\mathbb{E}[\hat{S}_{t+1}|S_t, A_t] = \mathbb{E}[\check{S}_{t+1}|S_t, A_t] \quad \forall (S_t, A_t) \in \mathcal{E}$ .*

$$\mathbb{E}[\hat{S}_{t+1}|S_t, A_t] = \mathbb{E}[\check{S}_{t+1}|S_t, A_t] \quad \forall (S_t, A_t) \in \mathcal{E}. \quad (10)$$

Figure 1 empirically shows that these assumptions are reasonable. The Infoprop distribution is slightly more stochastic than the ground truth process, which indicates that Assumption 1 holds. As (9) states, the model does not underestimate aleatoric uncertainty; the Infoprop rollouts should be at least as stochastic as the true process. Further, we observe no substantial bias of the Infoprop distribution underscoring the soundness of Assumption 2. Infoprop shows a similar behavior in high dimensional problems as reported in Section 6.

## 4.2 PREDICTIVE DISTRIBUTION

In the following, we aim to reduce the influence of stochasticity due to the uncertainty in  $\Theta_t$  on the predictive distribution of the model. Due to Eq. (3), we reformulate

$$\hat{S}_{t+1} = \mathbb{E}_{\hat{P}_{S,TS}}[S_{t+1}|S_t = s_t, A_t = a_t, W_t, N_t]$$

that attributes stochasticity in  $\hat{S}_{t+1}$  to aleatoric  $W_t$  and epistemic  $N_t$  noise. We shift the perspective from

## 4.2 DECOMPOSING THE MODEL STATE IN SIGNAL AND NOISE

We aim to isolate the stochasticity due to parameter uncertainty in  $\Theta_t$  to epistemic noise  $N_t$  as it allows us to analyze  $\hat{S}_{t+1}$ . We use the model error in  $\mathcal{S}$  rather than process (8) to project the noise in  $\vartheta$  and make the information-theoretic arguments that are at the core of.

Following (3) as well as Assumption 2, we approximate the distribution of  $\hat{S}_{t+1}$  by drawing  $\theta_t^e$ ,  $e \in \{1, \dots, E\}$  realizations from  $\Theta_t$  to infer an estimate of  $\check{S}_{t+1}$  to the same space as the signal, i.e. the dynamics, which is  $\mathcal{S}$ . The parameter distribution  $\Theta_t \sim \mathbb{P}_\Theta$  can induce arbitrarily complex distributions  $\Delta_t \sim \mathbb{P}_\Delta$ . To simplify the analysis, we solely consider the first two moments of  $\mathbb{P}_\Delta$ , namely  $\mu^\Delta$  and  $\Sigma^\Delta$ . We propose to improve the predictions by eliminating the epistemic uncertainty. Leveraging techniques from sensor fusion, we calculate This allows to reformulate the propagation equation (4) of the model state

$$\hat{S}_{t+1} = \check{S}_{t+1} + \Delta_t \approx \check{S}_{t+1} + \mu^\Delta(S_t, A_t) + L^\Delta(S_t, A_t)N_t \quad (11)$$

concerning the  $\check{S}_{t+1}$  and  $\Delta_t$  represented by  $\mu^\Delta(S_t, A_t)$  the model bias,  $\Sigma^\Delta(S_t, A_t)$  the epistemic variance with Cholesky decomposition  $L^\Delta(S_t, A_t)$ , and  $N_t$  the epistemic noise.

By Assumption 2, we have  $\mu^\Delta(S_t, A_t) = 0 \quad \forall (S_t, A_t) \in \mathcal{E}$ . Consequently, we can interpret the model rollout as communication through a Gaussian noise channel Cover & Thomas (2006) via (11).

Based on the propagation equation (11), we aim to infer the maximum likelihood estimate of  $\hat{S}_{t+1}$ , given the set of realizations  $\left\{ \mathbb{E} \left[ \hat{S}_{t+1} | \Theta_t = \theta_t^e \right] \right\}_{e=1}^E$ . Sampling  $\theta$  is straightforward and directly results in samples from  $\mathcal{N}$  by evaluating the neural network. Sampling directly from  $\mathcal{N}$  is intractable, however, access to  $\mathcal{N}$  allows us to use sensor fusion techniques, since the noise and signal are now in the same space. As the correlation between the realizations is unknown and significant we use Assumption 1 to obtain the predictive distribution using covariance intersection Julier & Uhlmann (2001).

$$\hat{S}_{t+1} = \mathbb{E} \left[ \hat{S}_{t+1} | N_t = 0 \right] \approx \bar{S}_{t+1} = \bar{\mu}(S_t, A_t) + \bar{L}(S_t, A_t) W_t$$

$E$  realizations of  $\left\{ \mathbb{E} \left[ \hat{S}_{t+1} | N_t = n_t^e \right] \right\}_{e=1}^E$ , to use it as the predictive distribution for our rollout scheme. As we cannot sample  $N_t$  directly, we instead use an equivalent definition of  $\hat{S}_{t+1}$ .

**Definition 4** (Model state concerning epistemic uncertainty). *Based on the model error process (8) the model state is defined as*

$$\hat{S}_{t+1} = \mathbb{E}_{\hat{P}_{S,TS}} [S_{t+1} | S_t = s_t, A_t = a_t, W_t, \Delta_t] \approx \mathbb{E}_{\hat{P}_{S,TS}} [S_{t+1} | S_t = s_t, A_t = a_t, W_t, N_t] \quad (12)$$

with  $\bar{L}(S_t, A_t)^\top \bar{L}(S_t, A_t) = \bar{\Sigma}(S_t, A_t)$

Reformulating (6) concerning  $\Delta_t$  does not change the information content or the induced sigma-algebra, as  $\Delta_t$  is a measurable function of  $\Theta_t$ . In the simplified setting of solely considering the first two moments of  $\mathbb{P}_\Delta$ ,  $\bar{\Sigma}(S_t, A_t) = \left( \frac{1}{E} \sum_{e=1}^E \hat{\Sigma}_{\Theta_t = \theta_t^e}(S_t, A_t)^{-1} \right)^{-1}$ , and  $\bar{\mu}(S_t, A_t) = \bar{\Sigma}(S_t, A_t) \left( \frac{1}{E} \sum_{e=1}^E \hat{\Sigma}_{\Theta_t = \theta_t^e}(S_t, A_t)^{-1} \hat{\mu}_{\Theta_t = \theta_t^e}(S_t, A_t) \right)$ . A derivation is provided in Appendix D.2.  $N_t$  fully describes stochasticity due to model error. In reverse, we can obtain realizations  $\left\{ \mathbb{E} \left[ \hat{S}_{t+1} | \Theta_t = \theta_t^e \right] \right\}_{e=1}^E$  and interpret them as samples  $\left\{ \mathbb{E} \left[ \hat{S}_{t+1} | N_t = n_t^e \right] \right\}_{e=1}^E$ .

**Lemma 1.** *Given  $E$  realizations of  $\mathbb{E} \left[ \hat{S}_{t+1} | \Theta_t = \theta_t^e \right]$ , we can estimate the environment state using maximum likelihood as*

$$\hat{S}_{t+1} = \mathbb{E} \left[ \hat{S}_{t+1} | N_t = 0 \right] \approx \bar{S}_{t+1} = \bar{\mu}(S_t, A_t) + \bar{L}(S_t, A_t) W_t \quad (13)$$

*Proof.* see Appendix D.2.1 □

**Lemma 2.** *Following this line of thought, the maximum likelihood estimate of  $\Sigma^\Delta$  is given by*

$$\bar{\Sigma}^\Delta(S_t, A_t) = \frac{1}{E} \sum_{e=1}^E \left( \hat{\mu}_{\Theta_t = \theta_t^e}(S_t, A_t) - \bar{\mu}(S_t, A_t) \right) \left( \hat{\mu}_{\Theta_t = \theta_t^e}(S_t, A_t) - \bar{\mu}(S_t, A_t) \right)^\top. \quad (14)$$

### 4.3 DATA CORRUPTION QUANTIFICATION

*Estimating-*  
*Proof.* see Appendix D.2.1 □

Given the maximum likelihood estimates of the environment state  $\bar{S}_{t+1}$  and the epistemic variance  $\bar{\Sigma}^\Delta$  makes (3) applicable to model-based rollouts, where we have no access to the true environment dynamics and model error. Intuitively, we perform a TS rollout, project the resulting state realizations  $\hat{s}_{t+1}$  into the maximum likelihood distribution  $\bar{S}_{t+1}$ , and estimate data corruption through model usage by the conditional entropy of, we can decompose the model state  $\hat{S}_{t+1}$  in a signal and noise fraction according to (11) in  $\mathcal{E}$ .



### 4.3 CONSTRUCTING THE INFOPROP STATE

Having decomposed  $\hat{S}_{t+1}$  into signal  $\bar{S}_{t+1}$  given the model realization  $\hat{s}_{t+1}$ . To do so we introduce and noise  $\tilde{\Sigma}^\Delta$ , allows us to define the Infoprop state.

**Definition 5** (Infoprop state). *We define the Infoprop state*

$$\tilde{S}_{t+1} := \mathbb{E} \left[ \bar{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1} \right] \equiv \mathbb{E}_{\tilde{P}_{S,IP}} \left[ S_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}, U_t \right] \quad (15)$$

as the conditional expectation of the *next-estimated environment* state given a *realization of the state-action pair and a model sample under the sample of the model state*. We derive the *corresponding* Infoprop kernel  $\tilde{P}_{S,IP}(\cdot | S_t, A_t, \hat{S}_{t+1}) = \mathcal{N} \left( \tilde{\mu}(S_t, A_t, \hat{S}_{t+1}), \tilde{\Sigma}(S_t, A_t, \hat{S}_{t+1}) \right)$  with the conditional noise  $U_t \sim \mathcal{N}(0, I)$ . *Following (3) and (3), an transition is defined by*

$$\tilde{S}_{t+1} = \tilde{\mu}(S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}) + \tilde{L}(S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1})U_t$$

with *mean*  $\tilde{\mu}(s_t, a_t, \hat{s}_{t+1}) = \bar{\mu}(s_t, a_t) + K(s_t, a_t) (\bar{L}(s_t, a_t)w_t + \bar{L}^\Delta(s_t, a_t)n_t)$ , *variance*  $\tilde{\Sigma}(s_t, a_t, \hat{s}_{t+1}) = (I - K(s_t, a_t))\tilde{\Sigma}(s_t, a_t)$ ,  $K(s_t, a_t) = \tilde{\Sigma}(s_t, a_t) (\tilde{\Sigma}(s_t, a_t) + \tilde{\Sigma}^\Delta(s_t, a_t))^{-1}$ ,  $\tilde{L}(s_t, a_t)\tilde{L}(s_t, a_t)^\top = \tilde{\Sigma}(s_t, a_t)$ , *and*  $\bar{L}^\Delta(s_t, a_t)\bar{L}^\Delta(s_t, a_t)^\top = \tilde{\Sigma}^\Delta(s_t, a_t)$ . *Figure 3 builds intuition to in Appendix D.3.*

Consequently, the Infoprop rollouts. Given a realization  $(s_t, a_t)$  and the parameter realizations  $\theta_t^e, e \in \{1, \dots, E\}$ , we approximate  $\hat{S}_{t+1}$  and compute state aims to infer the signal  $\bar{S}_{t+1}$  via (13). Sampling one of the parameter realizations  $\theta_t^{e'}$  and an aleatoric noise realization  $w_t$  provides us with a realization  $\hat{s}_{t+1}$  given a noisy observation  $\hat{s}_{t+1}$ . Propagating model-based rollouts using  $\tilde{S}_{t+1}$ , yields favorable properties as stated in Theorem 1.

**Theorem 1** (Infoprop state). *By construction,  $\tilde{S}_{t+1}$  addresses questions (i) and (ii) of Section 3.*

(i) *The distribution of Infoprop states is identical to the estimated environment distribution.*

$$\tilde{S}_{t+1} \stackrel{\text{dist}}{=} \bar{S}_{t+1} \quad (16)$$

*Proof.* see Appendix D.4.  $\square$

*Considering*

(ii) *The sum of marginal entropies of  $\tilde{S}_{t+1}$  defines the information loss along an Infoprop rollout.*

$$\mathbb{H}(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T) = \sum_{t=0}^T \mathbb{H}(\tilde{S}_{t+1}) \quad (17)$$

*Proof.* see Appendix D.5.  $\square$

Figure 3 illustrates the Infoprop rollout mechanism and provides intuition for Theorem 1. In the case of a perfect model, i.e.  $\Sigma^\Delta = 0$ , as  $\tilde{\Sigma}^\Delta = 0$ , depicted in Figure 3a, conditioning  $\tilde{S}_{t+1}$  on the realization  $\hat{s}_{t+1}$  essentially provides information about  $w_t$ , resulting in a deterministic state  $\tilde{s}_{t+1}$  as follows from (5). Thus,  $\mathbb{H}(\tilde{S}_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}) = \mathbb{H}(\tilde{S}_{t+1}) = 0$  in this scenario.

provides the information about the process noise realization  $w_t$  without ambiguity. Consequently, the belief about the environment state given the sample from the model  $\tilde{S}_{t+1} = \mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}]$  is a deterministic object and  $\mathbb{H}(\tilde{S}_{t+1}) = 0$ . In the general case of model error as scenario of  $\Sigma^\Delta > 0$  depicted in Figure 3b, the sample from the model  $\hat{s}_{t+1}$  leaves epistemic uncertainty results in ambiguity about the corresponding realization  $\tilde{s}_{t+1}$ , as the model error corrupts the information. Consequently,  $\mathbb{H}(\tilde{S}_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}) = \mathbb{H}(\tilde{S}_{t+1}) > 0$ , which corresponds to the information lost<sup>2</sup> due to using the faulty model for predicting the next state. Therefore, propagating environment state given  $\hat{s}_{t+1}$ , such that  $\mathbb{H}(\tilde{S}_{t+1}) > 0$ . Notably, conditioning  $\tilde{S}_{t+1}$  on  $\hat{s}_{t+1}$ , results in Infoprop predictions  $\tilde{S}_{t+1}$  as the rollout variable yields an estimate for data corruption due to model error.

Most importantly, this effect extends to the multi-step prediction setting depicted in Figure 3c. Rollouts are performed by propagating the following estimated environment distribution  $\tilde{S}_{t+1}$  as stated in Theorem 1 (i). This results in a data distribution that closely resembles the environment dynamics as desired in question (i) of Section 3. Finally, Figure 3c depicts a Infoprop state  $\tilde{s}_{t+1}$  obtained from a realization of conditional noise  $u_t$ . The information loss rollout propagated via realization  $\tilde{s}_{t+1}$ . We measure data corruption due to model error along a trajectory of length  $T$ .

$$\mathbb{H}(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T) = \sum_{t=0}^T \mathbb{H}(\tilde{S}_{t+1})$$

equals the sum of entropies  $\mathbb{H}(\tilde{S}_{t+1})$  along this path as derived Lemma D.5 of Appendix D.5 using the conditional entropy of a rollout under the estimated environment dynamics  $(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T)$  given the realizations observed from the model  $(S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T)$ , i.e. given the observed model trajectory, how sure are we on how the corresponding environment trajectory would look like?. As per Theorem 1 (ii), this trajectory-based approach to uncertainty can be addressed with the accumulated marginal entropy of  $\tilde{S}_{t+1}$ , addressing question (ii) of Section 3.

Consequently, propagating model-based rollouts with  $\tilde{P}_{\tilde{S}, \text{IP}}$  instead of the TS kernel  $\hat{P}_{\tilde{S}, \text{TS}}$  has two fundamental benefits. First, level sets of  $\tilde{S}_{t+1}$  lie within level sets of  $\bar{S}_{t+1}$  by construction, i.e. rollouts follow the improved predictive distribution derived in Section 4.2. Second, the accumulated entropy of the state corresponds to data corruption due to model error.

<sup>2</sup>In information-theoretic terms this rather corresponds to generated information but we believe this formulation is more intuitive from a MBRL practitioner’s perspective.

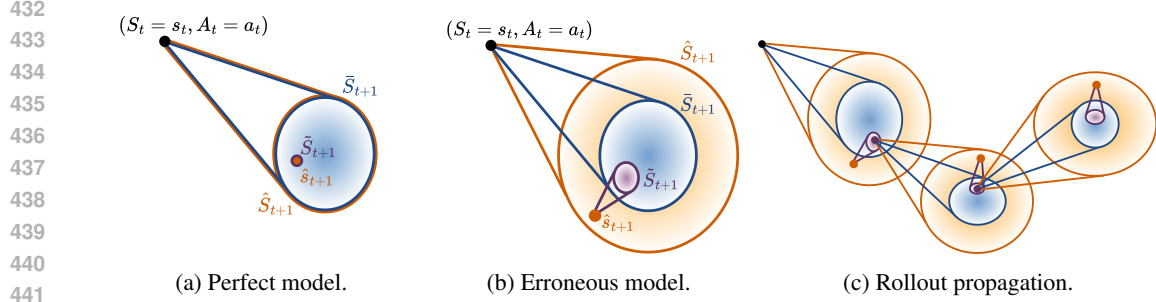


Figure 3: Infoprop rollout mechanism. (a), (b): Generating the Infoprop state  $\tilde{S}_{t+1}$  from the estimated predictive distribution  $\tilde{S}_{t+1}$  and the model sample  $\hat{s}_{t+1}$ . (c) Performing an Infoprop rollout.

#### 4.4 ROLLOUT TERMINATION CRITERIA

Having introduced how to propagate Infoprop rollouts, the question remains when to terminate them. In the following, we propose two termination criteria to address question (iii) of Section 3.

First, Infoprop rollouts build on the assumption that ~~model-based rollouts are performed in model usage is restricted to~~ a sufficiently accurate subset  $\mathcal{E} \subseteq \mathcal{S} \times \mathcal{A}$ . ~~Following, following~~ the ideas of Frauenknecht et al. (2024), ~~we construct  $\mathcal{E}$  on the notion of single-step predictive uncertainty.~~ ~~Consequently, we define  $\mathcal{E} := \{(s_t, a_t) \in \mathcal{S} \times \mathcal{A} \mid \mathbb{H}(\tilde{S}_{t+1}) \leq \lambda_1, \hat{s}_{t+1} \sim \hat{P}_{S,TS}(\cdot | s_t, a_t)\}$  amenable to the setting, with  $\lambda_1$ .~~

**Definition 6** (Sufficiently accurate subset). We define the threshold of sufficiently accurate subset

$$\mathcal{E} := \{(s_t, a_t) \in \mathcal{S} \times \mathcal{A} \mid \mathbb{H}(\tilde{S}_{t+1}) \leq \lambda_1, \hat{s}_{t+1} \sim \hat{P}_{S,TS}(\cdot | s_t, a_t)\} \quad (18)$$

based on a threshold  $\lambda_1$  for the single-step information loss represented by  $\mathbb{H}(\tilde{S}_{t+1})$ . ~~Therefore, Second, we restrict~~ Infoprop rollouts ~~are terminated, whenever  $\mathbb{H}(\tilde{S}_{t+1}) > \lambda_1$  to enforce operation in  $\mathcal{E}$ .~~

~~Second, accumulating smaller errors can corrupt the data distribution of long model rollouts. Thus to sufficiently accurate paths to limit uncertainty accumulation.~~

**Definition 7** (Sufficiently accurate path). Based on the estimated information loss along a rollout (4.3), we define an upper bound on cumulative entropy  $\sum_t \mathbb{H}(\tilde{S}_{t+1}) \leq \lambda_2$  to avoid data corruption. the set of sufficiently accurate paths of length  $t' \in \{1, \dots, T\}$  as

$$\mathcal{P}^{t'} := \left\{ (s_t, a_t)_{t=0}^{t'} \in (\mathcal{S} \times \mathcal{A})^{t'} \mid \sum_{t=0}^{t'} \mathbb{H}(\tilde{S}_{t+1}) \leq \lambda_2 \right\}. \quad (19)$$

Heuristics for determining values of  $\lambda_1$  and  $\lambda_2$  depend on the class of AES model and MBRL algorithm at hand. ~~An example is with an example~~ provided in Section 5. Combining the steps above yields the Infoprop rollout mechanism illustrated in Algorithm 1.

## 5 AUGMENTING STATE-OF-THE-ART: INFOPROP-DYNA

While the Infoprop rollout mechanism is applicable to different kinds of MBRL with AES models, we illustrate its capabilities in a Dyna-style architecture with probabilistic ensemble (PE) models Lakshminarayanan et al. (2017). We design *Infoprop-Dyna* by integrating the Infoprop rollout mechanism in the state-of-the-art framework proposed in Janner et al. (2019) with minor adaptations.

As discussed in Section 4.4, heuristics for  $\lambda_1$  and  $\lambda_2$  ~~need to be designed for~~ depend on the algorithm at hand. In Infoprop-Dyna, we take the common approach Chua et al. (2018); Janner et al. (2019) of neglecting cross-correlations between state dimensions for computational reasons. Thus, ~~all covariances are diagonal matrices and~~ we can consider the data corruption of each state dimension independently. As the predictive quality of different state dimensions can differ substantially,

**Algorithm 1** Infoprop

---

```

486
487
488 Require:  $s_0$ 
489 while  $t < T + 1$  do
490    $a_t \sim \pi(\cdot | s_t)$ 
491   for  $e \in \{1, \dots, E\}$  do
492      $\theta_t^e \sim \mathbb{P}_\Theta$ 
493      $\bar{S}_{t+1}(s_t, a_t)$  from (13), and  $\bar{\Sigma}^\Delta(s_t, a_t)$  from (14)
494      $\hat{s}_{t+1} = \mathbb{E} \left[ \hat{S}_{t+1} | W_t = w_t, \Theta_t = \theta_t^e \right]$  with  $w_t \sim \mathcal{N}(0, I)$ ,  $\theta_t^e \sim \mathcal{U}(\{\theta_t^1, \dots, \theta_t^E\})$ 
495      $\tilde{S}_{t+1}$  from (5) and  $\mathbb{H}(\tilde{S}_{t+1})$  from (24)
496     if  $\mathbb{H}(\tilde{S}_{t+1}) > \lambda_1$  then
497       break
498     else if  $\sum_{t'=0}^t \mathbb{H}(\tilde{S}_{t'+1}) > \lambda_2$  then
499       break
500     else
501        $s_t \leftarrow \mathbb{E}[\tilde{S}_{t+1} | U_t = u_t]$  with  $u_t \sim \mathcal{N}(0, I)$ 

```

---

we choose both thresholds to be as  $n_S$  dimensional vectors, such that a rollout is terminated as soon as the data corruption of any state-dimension overshoots the corresponding threshold.

In Dyna-style MBRL Janner et al. (2019), the dynamics model is trained on the data distribution observed during environment interaction. The corresponding transitions are stored in an environment replay buffer  $\mathcal{D}_{\text{env}} = \left\{ \left( \check{s}_t^{(b)}, \check{a}_t^{(b)}, \check{r}_{t+1}^{(b)}, \check{s}_{t+1}^{(b)} \right) \right\}_{b=1}^{|\mathcal{D}_{\text{env}}|}$   $\mathcal{D}_{\text{env}} = \left\{ \left( \check{s}_t^{(b)}, \check{a}_t^{(b)}, \check{r}_{t+1}^{(b)}, \check{s}_{t+1}^{(b)} \right) \right\}_{b=1}^{|\mathcal{D}_{\text{env}}|}$ , where  $(b)$  indicates the index in the replay buffer. After a fixed number of interaction steps between a model-free RL agent and the environment, the dynamics model is retrained on the data in  $\mathcal{D}_{\text{env}}$ , model-based rollouts are performed, and the data is stored in a replay buffer  $\mathcal{D}_{\text{mod}}$  to train the model-free RL agent. Consequently, we assume the PE model to be accurate within the data distribution of  $\mathcal{D}_{\text{env}}$  and build the heuristic for  $\lambda_1$  and  $\lambda_2$  on the predictive uncertainty within the environment buffer.

After each round of retraining the PE model, we compute a set of dimension-wise Infoprop state entropies for single-step predictions in  $\mathcal{D}_{\text{env}}$  according to

$$\mathcal{H}^k = \left\{ \mathbb{H} \left( \bar{S}_{t+1}^k | S_t = \check{s}_t^{(b)}, A_t = \check{a}_t^{(b)}, \hat{S}_{t+1}^k = \hat{s}_{t+1}^{k,(b)} \right) \right\}_{b=1}^{|\mathcal{D}_{\text{env}}|} \quad (20)$$

where  $k \in \{1, \dots, n_S\}$  indicates the corresponding state dimension. We define the dimension-wise thresholds  $\lambda_1^k$  and  $\lambda_2^k$  based on the cumulative distribution function of dimension-wise entropies

$$F_{\mathcal{H}^k}(h) = \frac{1}{|\mathcal{H}^k|} \sum_{h' \in \mathcal{H}^k} \mathbb{1}[h' \leq h]. \quad (21)$$

We define the  $k^{\text{th}}$  element of  $\lambda_1$  is defined as the  $\zeta_1$  quantile of the single-step entropy set

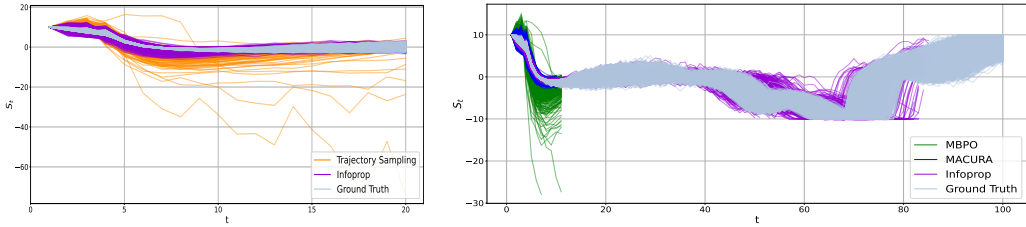
$$\lambda_1^k = \inf \{ h \in \mathcal{H}^k : F_{\mathcal{H}^k}(h) \geq \zeta_1 \} \quad (22)$$

to limit and limits model usage to the sufficiently accurate subset  $\mathcal{E}$ . Similarly, To restrict rollouts of length  $t'$  to  $\mathcal{P}^{t'}$ , we define the  $k^{\text{th}}$  element of  $\lambda_2$  is defined by as the  $\zeta_2$  quantile of the entropy set scaled by  $\xi$

$$\lambda_2^k = \xi \inf \{ h \in \mathcal{H}^k : F_{\mathcal{H}^k}(h) \geq \zeta_2 \}. \quad (23)$$

Here,  $\zeta_2$  denotes a quantile corresponding to precise predictions and  $\xi$  to the number of prediction steps we are willing to accumulate the resulting data corruption. We choose  $\zeta_1 = 0.99$ ,  $\zeta_2 = 0.01$  and  $\xi = 100$  for all experiments in Section 6 without further hyperparameter tuning.

We use pink noise for environment exploration Eberhard et al. (2023) to quickly expand  $\mathcal{E}$  Frauenknecht et al. (2024). Pseudocode is provided in Algorithm 3 of Appendix C.



(a) Trajectory Sampling vs. Infoprop

(b) MBPO vs. MACURA vs. Infoprop-Dyna

Figure 4: Predictive quality of rollouts in the 11<sup>th</sup> state dimension of MuJoCo hopper. (a) Rollouts according to Trajectory Sampling (TS) and Infoprop. (b) Rollout schemes of MBPO and MACURA based on TS compared to Infoprop-Dyna.

## 6 EXPERIMENTS AND DISCUSSION

To demonstrate the benefits of the Infoprop mechanism, we compare Infoprop-Dyna to state-of-the-art Dyna-style MBRL algorithms on MuJoCo Todorov et al. (2012) benchmark tasks. We report

- substantial improvements in the consistency of predicted data, especially over long horizons;
- effective rollout termination based on accumulated model error propagation; and
- state-of-the-art performance in Dyna-style MBRL on several MuJoCo tasks.

Furthermore, we discuss the limitations of naively integrating Infoprop into the standard Dyna-style setup Janner et al. (2019) and point to further research questions.

### 6.1 EXPERIMENTAL SETUP

We compare Infoprop-Dyna to Model-Based Policy Optimization (MBPO) Janner et al. (2019) and Model-Based Actor-Critic with Uncertainty-Aware Rollout Adaption (MACURA) Frauenknecht et al. (2024) as well as to Soft Actor-Critic (SAC) Haarnoja et al. (2018) that represents the model-free learner of all the Dyna-style approaches above. We build our implementation<sup>2</sup> on the code base<sup>3</sup> provided by Frauenknecht et al. (2024). Further details are provided in Appendix E.1

### 6.2 PREDICTION QUALITY

To compare different rollout mechanisms, we train an Infoprop-Dyna agent on hopper for 120000 environment interactions and perform model rollouts from states in  $\mathcal{D}_{env}$ .

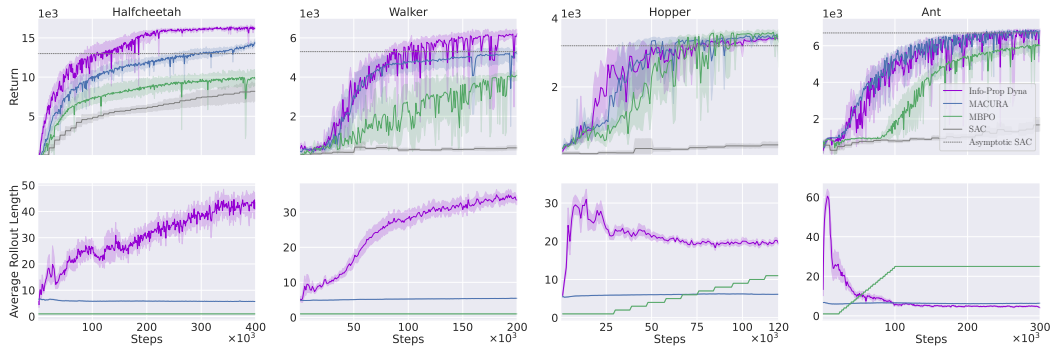
First, we evaluate the consistency of Infoprop and TS rollouts, propagating 20 steps without termination. Figure 4a depicts the resulting distributions for the 11<sup>th</sup> dimension of the hopper state. Infoprop rollouts follow the ground truth distribution closely and show substantially improved data consistency compared to TS rollouts. This underscores the improved predictive distribution, underscoring the ability of Infoprop that effectively mitigates to effectively mitigate model error propagation.

Next, we compare the rollout mechanisms of MBPO and MACURA based on TS sampling with Infoprop-Dyna rollouts. Figure 4b shows the results for 11<sup>th</sup> dimension of the hopper and a maximum rollout length of 100 steps. MBPO rollouts are propagated for 11 steps following the schedule proposed in Janner et al. (2019), resulting in a widely spread distribution. In contrast, MACURA has an adaptive rollout length capped at 10 steps Frauenknecht et al. (2024), leading to better data consistency. The improved predictive distribution and capability to estimate accumulated error of Infoprop allows for substantially longer rollouts up to 100 steps. The Infoprop termination criteria reliably stop distorted rollouts, resulting in consistent rollouts over long horizons. Appendix E.2 provides additional results for setting the maximum rollout length of all three approaches to 100.

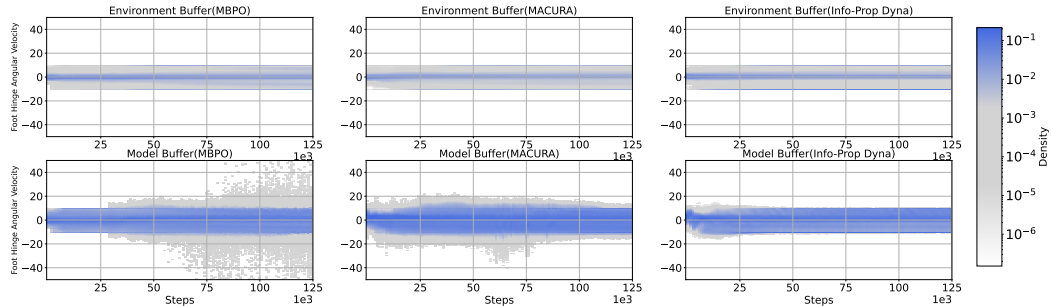
<sup>2</sup>Code will be published upon acceptance and is currently provided in the supplementary material.

<sup>3</sup><https://github.com/Data-Science-in-Mechanical-Engineering/macura>

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647



(a) Performance and average rollout length on MuJoCo tasks.



(b) Adequacy of  $\mathcal{D}_{\text{mod}}$  on the 11<sup>th</sup> state dimension of hopper.

Figure 5: Evaluation on MuJoCo tasks. (a) *Infoprop-Dyna* shows state-of-the-art performance for Dyna-style MBRL on several MuJoCo tasks while considerably increasing average rollout length on most tasks. (b) *Infoprop-Dyna* shows substantially improved consistency between  $\mathcal{D}_{\text{env}}$  and  $\mathcal{D}_{\text{mod}}$ .

### 6.3 PERFORMANCE EVALUATION

As depicted in the top row of Figure 5a, Infoprop-Dyna performs on par with or better than MACURA, while substantially outperforming MBPO with respect to data efficiency and asymptotic performance. Notably, Infoprop-Dyna consistently outperforms SAC with a fraction of environment interaction. The bottom row of Figure 5a depicts the average rollout lengths. Infoprop-Dyna shows substantially increased rollout lengths compared to prior methods in all environments but ant.

A major concern of this work is the consistency of model-based rollouts with the environment distribution. Figure 5b depicts the data distribution in  $\mathcal{D}_{\text{env}}$  and  $\mathcal{D}_{\text{mod}}$  of the respective Dyna-style approaches throughout training for the 11<sup>th</sup> dimension of the hopper state. The figure shows a histogram over state values over the course of training. It can be seen that the model data distribution of Infoprop-Dyna closely follows the distribution observed in the environment, while both the data from MBPO and MACURA show severe outliers. This is the case, even though the rollout data in Infoprop-Dyna is obtained from substantially longer rollouts as can be seen from Figure 5a which indicates the capabilities of the Infoprop rollout mechanism.

### 6.4 LIMITATIONS AND OUTLOOK

Despite the excellent quality of model-generated data with the Infoprop rollout, the limitations of Infoprop-Dyna are most apparent on MuJoCo humanoid with results provided in E.3. These show instabilities in learning and point to structural problems when integrating Infoprop rollouts naively into standard Dyna-style architectures Janner et al. (2019).

Figure 5b shows that the long rollouts of Infoprop-Dyna can cause rapid distribution shifts in  $\mathcal{D}_{\text{mod}}$ , especially early in training. These nonstationary buffers are a challenge to deep Q-learning methods Mnih et al. (2015). **The main issue with Infoprop-Dyna is likely overfitting critics and plasticity loss** Nikishin et al. (2022); D’Oro et al. (2023), as also reported by Frauenknecht et al. (2024) for Dyna-style MBRL trained on high-quality data. **Another challenge** **Another issue** is primacy bias in model learning Qiao et al. (2023), where the model overfits to initial data and subsequently struggles to generalize, as seen in the decreasing rollout length for the ant environment

in Figure 5a. [The main problem with Infoprop-Dyna is likely overfitting critics and plasticity loss Nikishin et al. \(2022\); D’Oro et al. \(2023\), as also reported by Frauenknecht et al. \(2024\) for Dyna-style MBRL trained on high-quality data. We provide an ablation on this observation and sketch methods to counteract this phenomenon in Appendix E.4.](#)

## 7 RELATED WORK

The negative effects of accumulated model error on the performance of MBRL methods is a long-studied problem Venkatraman et al. (2015); Talvitie (2016); Asadi et al. (2018b;a).

Different model architectures have been proposed to mitigate this issue, such as trajectory models Asadi et al. (2019); Lambert et al. (2021), bidirectional models Lai et al. (2020), temporal segment models Mishra et al. (2017) or self-correcting models Talvitie (2016). These architectures, however, imply substantial additional effort for model learning, such that state-of-the-art performance in the respective fields of MBRL is often reported for simpler single-step model architectures Chua et al. (2018); Janner et al. (2019); Buckman et al. (2018).

These approaches address the problem of error accumulation by keeping model-based rollouts sufficiently short. Janner et al. (2019) introduce the concept of branched rollouts that allows to cover relevant parts of  $\mathcal{S}$  with short model rollouts. Other methods weight rollouts of different lengths according to their single-step uncertainty Buckman et al. (2018) or use single-step uncertainty to schedule rollout length Pan et al. (2020); Frauenknecht et al. (2024). Infoprop allows to infer model data consistent with the environment distribution over long rollout horizons using comparatively simple model architectures and computationally cheap conditioning operations.

Infoprop is inspired by an information-theoretic view on RL Lu et al. (2023). Thus far, information-theoretic arguments have been mostly used to improve the exploration Haarnoja et al. (2018); Lu & Roy (2019); Ahmed et al. (2019); Mohamed & Rezende (2015) and generalization Tishby & Zaslavsky (2015); Lu et al. (2020); Igl et al. (2019); Islam et al. (2023) of model-free RL methods. While aspects of dynamical systems such as causality, modeling, and control Lozano-Duran & Aranz (2021), predictability Kleeman (2011) or dealing with noisy observations Gattami (2014) have been studied from an information theoretic perspective, these works do not directly apply to the MBRL setup nor extend to long model-based rollouts.

## REFERENCES

- Zafarali Ahmed, Nicolas Le Roux, Mohammad Norouzi, and Dale Schuurmans. Understanding the impact of entropy on policy optimization. *Int. Conf. on Machine Learning*, abs/1811.11214, 2019.
- Kavosh Asadi, Evan Cater, Dipendra Misra, and Michael L. Littman. Towards a Simple Approach to Multi-step Model-based Reinforcement Learning. *arXiv*, October 2018a. doi: 10.48550/arXiv.1811.00128.
- Kavosh Asadi, Dipendra Misra, and Michael L. Littman. Lipschitz Continuity in Model-based Reinforcement Learning. *arXiv*, April 2018b. doi: 10.48550/arXiv.1804.07193.
- Kavosh Asadi, Dipendra Misra, Seungchan Kim, and Michel L. Littman. Combating the Compounding-Error Problem with a Multi-step Model. *arXiv*, May 2019. doi: 10.48550/arXiv.1905.13320.
- Philipp Becker and Gerhard Neumann. On Uncertainty in Deep State Space Models for Model-Based Reinforcement Learning. *Transactions on Machine Learning Research*, October 2022. doi: 10.48550/arXiv.2210.09256.
- Marc G. Bellemare, Will Dabney, and Mark Rowland. *Distributional Reinforcement Learning*. MIT Press, 2023.
- Thomas Bi and Raffaello D’Andrea. Sample-efficient learning to solve a real-world labyrinth game using data-augmented model-based reinforcement learning. In *IEEE Int. Conf. on Robotics and Automation*, pp. 7455–7460, 2024. doi: 10.1109/ICRA57147.2024.10610577.

- 702 Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-  
703 efficient reinforcement learning with stochastic ensemble value expansion. In *Int. Conf. on Neural*  
704 *Information Processing Systems*. 2018.
- 705 Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep Reinforcement  
706 Learning in a Handful of Trials using Probabilistic Dynamics Models. *Adv. in Neural Information*  
707 *Processing Systems*, 2018.
- 708 Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory 2nd Edition (Wiley Series in*  
709 *Telecommunications and Signal Processing) by Thomas M. Cover Joy A. Thomas(2006-07-18)*.  
710 Wiley-Interscience, Hoboken, NJ, USA, January 2006.
- 711 Marc Peter Deisenroth and Carl Edward Rasmussen. PILCO: a model-based and data-efficient  
712 approach to policy search. In *Int. Conf. on Machine Learning*. 2011.
- 713 Pierluca D’Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G. Bellemare, and  
714 Aaron C. Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier.  
715 In *Int. Conf. on Learning Representations*, 2023.
- 716 Onno Eberhard, Jakob Hollenstein, Cristina Pinneri, and Georg Martius. Pink noise is all you need:  
717 Colored noise exploration in deep reinforcement learning. In *Int. Conf. on Learning Representa-*  
718 *tions*, 2023.
- 719 Vladimir Feinberg, Alvin Wan, Ion Stoica, Michael I. Jordan, Joseph E. Gonzalez, and Sergey  
720 Levine. Model-Based Value Estimation for Efficient Model-Free Reinforcement Learning. *Int.*  
721 *Conf. on Machine Learning*, 2018.
- 722 Bernd Frauenknecht, Artur Eisele, Devdutt Subhasish, Friedrich Solowjow, and Sebastian Trimpe.  
723 Trust the Model Where It Trusts Itself – Model-Based Actor-Critic with Uncertainty-Aware Roll-  
724 out Adaption. *Int. Conf. on Machine Learning*, May 2024. doi: 10.48550/arXiv.2405.19014.
- 725 Ather Gattami. Kalman meets shannon. *ArXiv*, 2014.
- 726 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash  
727 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic Algo-  
728 rithms and Applications. *arXiv*, 2018.
- 729 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James  
730 Davidson. Learning Latent Dynamics for Planning from Pixels. In *Int. Conf. on Machine Learn-*  
731 *ing*. PMLR, May 2019.
- 732 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to Control: Learn-  
733 ing Behaviors by Latent Imagination. In *Int. Conf. on Learning Representations*, April 2020.
- 734 Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari with  
735 Discrete World Models. *Int. Conf. on Learning Representations*, October 2021.
- 736 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering Diverse Domains  
737 through World Models. *Int. Conf. on Learning Representations*, January 2023.
- 738 Maximilian Igl, Kamil Ciosek, Yingzhen Li, Sebastian Tschitschek, Cheng Zhang, Sam Devlin,  
739 and Katja Hofmann. Generalization in reinforcement learning with selective noise injection and  
740 information bottleneck. In *Int. Conf. on Neural Information Processing Systems*, 2019.
- 741 Riashat Islam, Hongyu Zang, Manan Tomar, Aniket Didolkar, Md Mofijul Islam, Samin Yeasar  
742 Arnob, Tariq Iqbal, Xin Li, Anirudh Goyal, Nicolas Heess, and Alex Lamb. Representation  
743 learning in deep rl via discrete information bottleneck. In *Int. Conf. on Artificial Intelligence and*  
744 *Statistics*, 2023.
- 745 Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: model-  
746 based policy optimization. In *Int. Conf. on Neural Information Processing Systems*. 2019.
- 747 Simon Julier and Jeffrey Uhlmann. *General Decentralized Data Fusion with Covariance Intersec-*  
748 *tion (CI)*. 06 2001. doi: 10.1201/9781420038545.ch12.



- 756 Elia Kaufmann, Leonard Bauersfeld, Antonio Loquercio, Matthias Müller, Vladlen Koltun, and  
757 Davide Scaramuzza. Champion-level drone racing using deep reinforcement learning. *Nature*,  
758 August 2023.
- 759 Richard Kleeman. Information theory and dynamical system predictability. *Entropy*, 2011. doi:  
760 10.3390/e13030612.
- 761 Hang Lai, Jian Shen, Weinan Zhang, and Yong Yu. Bidirectional Model-based Policy Optimization.  
762 In *Int. Conf. on Machine Learning*. PMLR, 2020.
- 763 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
764 uncertainty estimation using deep ensembles. In *Int. Conf. on Neural Information Processing*  
765 *Systems*. 2017.
- 766 Nathan O. Lambert, Albert Wilcox, Howard Zhang, Kristofer S. J. Pister, and Roberto Calandra.  
767 Learning Accurate Long-term Dynamics for Model-based Reinforcement Learning. *IEEE Conf*  
768 *on Decision and Control*, December 2021. doi: 10.48550/arXiv.2012.09156.
- 769 Adrian Lozano-Duran and Gonzalo Arranz. Information-theoretic formulation of dynamical sys-  
770 tems: causality, modeling, and control. *ArXiv*, 2021.
- 771 Xingyu Lu, Kimin Lee, P. Abbeel, and Stas Tiomkin. Dynamics generalization via information  
772 bottleneck in deep reinforcement learning. *ArXiv*, 2020.
- 773 Xiuyuan Lu and Benjamin Van Roy. Information-theoretic confidence bounds for reinforcement  
774 learning. *Int. Conf. on Neural Information Processing Systems*, 2019.
- 775 Xiuyuan Lu, Benjamin Van Roy, Vikranth Dwaracherla, Morteza Ibrahimi, Ian Osband, and Zheng  
776 Wen. Reinforcement learning, bit by bit. *Foundations and Trends® in Machine Learning*, 2023.  
777 doi: 10.1561/22000000097.
- 778 Carlos E. Luis, Alessandro G. Bottero, Julia Vinogradska, Felix Berkenkamp, and Jan Peters.  
779 Model-Based Epistemic Variance of Values for Risk-Aware Policy Optimization. *arXiv*, 2023.
- 780 Nikhil Mishra, Pieter Abbeel, and Igor Mordatch. Prediction and Control with Temporal Segment  
781 Models. *Int. Conf. on Machine Learning*, March 2017. doi: 10.48550/arXiv.1703.04070.
- 782 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-  
783 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,  
784 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan Wier-  
785 stra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.  
786 *Nature*, 2015.
- 787 Shakir Mohamed and Danilo J. Rezende. Variational information maximisation for intrinsically  
788 motivated reinforcement learning. In *Int. Conf. on Neural Information Processing Systems*, 2015.
- 789 Anusha Nagabandi, Gregory Kahn, Ronald S. Fearing, and Sergey Levine. Neural Network Dynam-  
790 ics for Model-Based Deep Reinforcement Learning with Model-Free Fine-Tuning. In *IEEE Int.*  
791 *Conf. on Robotics and Automation*. 2018.
- 792 Michal Nauman, Michał Bortkiewicz, Piotr Miłoś, Tomasz Trzcinski, Mateusz Ostaszewski, and  
793 Marek Cygan. Overestimation, Overfitting, and Plasticity in Actor-Critic: the Bitter Lesson of  
794 Reinforcement Learning. In *International Conference on Machine Learning*, pp. 37342–37364.  
795 PMLR, July 2024. URL [https://proceedings.mlr.press/v235/nauman24a.](https://proceedings.mlr.press/v235/nauman24a.html)  
796 [html](https://proceedings.mlr.press/v235/nauman24a.html).
- 797 Evgenii Nikishin, Max Schwarzer, Pierluca D’Oro, Pierre-Luc Bacon, and Aaron C. Courville. The  
798 primacy bias in deep reinforcement learning. In *Int. Conf. on Machine Learning*, 2022.
- 799 OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw De-  
800 biak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal  
801 Jozefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé  
802 de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon  
803 Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep  
804 reinforcement learning. 2019.

810 Feiyang Pan, Jia He, Dandan Tu, and Qing He. Trust the model when it is confident: masked  
811 model-based actor-critic. In *Int. Conf. on Neural Information Processing Systems*. 2020.  
812

813 Zhongjian Qiao, Jiafei Lyu, and Xiu Li. Mind the Model, Not the Agent: The Primacy Bias in  
814 Model-based RL. *arXiv*, October 2023. doi: 10.48550/arXiv.2310.15017.

815 C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, July  
816 1948.

817

818 Dan Simon. *Optimal State Estimation*. January 2006. ISBN 978-0-47170858-2. doi: 10.1002/  
819 0470045345.

820

821 Laura Smith, Ilya Kostrikov, and Sergey Levine. Demonstrating A Walk in the Park: Learning to  
822 Walk in 20 Minutes With Model-Free Reinforcement Learning. *Robotics, Science and Systems*,  
823 August 2023. doi: 10.48550/arXiv.2208.07860.

824

825 Richard S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *SIGART  
826 Bull.*, 1991.

827

828 Erik Talvitie. Self-Correcting Models for Model-Based Reinforcement Learning. *arXiv*, December  
829 2016. doi: 10.48550/arXiv.1612.06018.

830

831 Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. *IEEE  
832 Information Theory Workshop (ITW)*, 2015.

833

834 Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based control.  
835 In *Int. Conf. on Intelligent Robots and Systems*. IEEE, 2012.

836

837 Miguel Vasco, Takuma Seno, Kenta Kawamoto, Kaushik Subramanian, Peter R. Wurman, and Peter  
838 Stone. A Super-human Vision-based Reinforcement Learning Agent for Autonomous Racing in  
839 Gran Turismo. *Reinforcement Learning Conference*, June 2024. doi: 10.48550/arXiv.2406.12563.

840

841 Arun Venkatraman, Martial Hebert, and J.. Bagnell. Improving Multi-Step Prediction of Learned  
842 Time Series Models. *AAAI*, February 2015. doi: 10.1609/aaai.v29i1.9590.

843

844 Grady Williams, Nolan Wagener, Brian Goldfain, Paul Drews, James M. Rehg, Byron Boots, and  
845 Evangelos A. Theodorou. Information theoretic MPC for model-based reinforcement learning. In  
846 *IEEE Int. Conf. on Robotics and Automation*. IEEE, 2017.

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

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

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

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  
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

## A NOTATION

### A.1 OBJECTS

- $S_t$  Random variable of a general state
- $S_t$  Random variable of the environment state
- $S_t$  Random variable of the estimated environment state
- $S_t$  Random variable of the model state
- $A_t$  Random variable of the action

... ...  
we will finish this for a potential camera-ready version.

## B TOY EXAMPLE

In Figure 1, we illustrate the data consistency of Trajectory Sampling Chua et al. (2018) and Infoprop in a one-dimensional random walk example with  $\mathcal{S} \subseteq \mathbb{R}$  and  $\mathcal{A} \subseteq \mathbb{R}$ . The dynamics follow (3) with  $\mu(S_t, A_t) = S_t + A_t$  and  $L(S_t, A_t) = 0.01$ . Actions are distributed according to  $A_t \sim \mathcal{N}(0, 0.1)$ . All rollouts start from  $s_0 = 0$  and are propagated for 100 steps. We perform 1000 rollouts under the environment dynamics and train a Probabilistic Ensemble Lakshminarayanan et al. (2017) model according to the information provided in Table 1. Subsequently, we perform 1000 model-based rollouts with this model and the respective rollout mechanism.

| Hyperparameter             | Value   |
|----------------------------|---------|
| number of ensemble members | 5       |
| number of hidden neurons   | 2       |
| number of layers           | 1       |
| learning rate              | 0.001   |
| weight decay               | 0.00001 |
| number of epochs           | 4       |

Table 1: Hyperparameters used for training the model on the random walk dataset.

## C PSEUDOCODE ALGORITHMS

---

### Algorithm 2 Trajectory Sampling Chua et al. (2018)

---

**Require:**  $s_0$   
**while**  $t < T + 1$  **do**  
 $a_t \sim \pi(\cdot | s_t)$   
 $\hat{s}_{t+1} = \mathbb{E} \left[ \hat{\mathcal{S}}_{t+1} | W_t = w_t, \Theta_t = \theta_t \right]$  with  $w_t \sim \mathcal{N}(0, I)$  and  $\theta_t \sim \mathbb{P}_\Theta$   
 $s_t \leftarrow \hat{s}_{t+1}$

---



---

### Algorithm 3 Infoprop-Dyna (Pseudocode adapted from Janner et al. (2019))

---

**Require:** Policy  $\pi$ , predictive AES model  $p_\Theta$ , environment buffer  $\mathcal{D}_{\text{env}}$ , model buffer  $\mathcal{D}_{\text{mod}}$ , rollout parameters  $T, \zeta_1, \zeta_2, \xi$   
**for**  $N$  epochs **do**  
  **for**  $J$  steps **do**  
    Interact with the environment according to  $\pi$ ; add to  $\mathcal{D}_{\text{env}}$   
    Train model  $p_\Theta$  on  $\mathcal{D}_{\text{env}}$   
    Perform single-step predictions with  $p_\Theta$  in  $\mathcal{D}_{\text{env}}$   
    Compute  $\lambda_1$  (22) and  $\lambda_2$  (23)  
    **for**  $M$  model rollouts **do**  
      Sample  $s_0$  uniformly from  $\mathcal{D}_{\text{env}}$   
      Perform Infoprop rollouts according to Algorithm 1; add to  $\mathcal{D}_{\text{mod}}$   
    **for**  $G \cdot J$  gradient updates **do**  
      Update  $\pi$  on  $\mathcal{D}_{\text{env}} \cup \mathcal{D}_{\text{mod}}$

---

## 972 D DERIVATIONS

### 973 D.1 QUANTIZED ENTROPY

974 For a RV  $Z \in \mathcal{Z} \subset \mathbb{R}^{n_z}$  with  $Z \sim \mathcal{N}(\mu_Z, \Sigma_Z)$  and discretization step size  $\Delta_z^{(k)}$  of the  $k^{\text{th}}$   
 975 dimension, the quantized entropy Cover & Thomas (2006) is

$$976 \mathbb{H}(Z) = \frac{1}{2} \log_2 ((2\pi e)^{n_z} |\Sigma_Z|) - \sum_{k=1}^{n_z} \log_2 (\Delta_z^{(k)}). \quad (24)$$

### 982 D.2 MAXIMUM LIKELIHOOD PREDICTIVE DISTRIBUTION

#### 983 D.2.1 PROOF OF LEMMA 1

984 *Proof.* We introduce the conditional expectation over the next state under the model, given a realization  $\theta_t^e$

$$985 \hat{S}_{t+1}^e := \mathbb{E}_{\hat{P}_{S,TS}} [\hat{S}_{t+1} | \Theta_t = \theta_t^e]. \quad (25)$$

986 Further,  $\hat{\mu}^e := \hat{\mu}_{\Theta_t = \theta_t^e}$ ,  $\hat{\Sigma}^e := \hat{\Sigma}_{\Theta_t = \theta_t^e}$  and  $\hat{L}^e := \hat{L}_{\Theta_t = \theta_t^e}$  such that

$$987 \hat{S}_{t+1}^e = \hat{\mu}^e(S_t, A_t) + \hat{L}^e(S_t, A_t)W_t. \quad (26)$$

988 Given  $E$  RVs  $\hat{S}_{t+1}^e$  we define their joint distribution

$$989 \begin{aligned} & \begin{pmatrix} \hat{S}_{t+1}^1 \\ \vdots \\ \hat{S}_{t+1}^E \end{pmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \hat{\mu}^1 \\ \vdots \\ \hat{\mu}^E \end{bmatrix}, \begin{bmatrix} \hat{\Sigma}^1 & \dots & \hat{\Sigma}^{1E} \\ \vdots & \ddots & \vdots \\ \hat{\Sigma}^{E1} & \dots & \hat{\Sigma}^E \end{bmatrix} \right) \\ & =: \hat{S} \sim \mathcal{N}(\hat{\mu}, \hat{\Sigma}) \end{aligned} \quad (27)$$

1000 with  $\hat{\Sigma}^{ef} := \text{Cov}[\hat{S}_{t+1}^e, \hat{S}_{t+1}^f]$ . We aim to track  $S_{t+1}$  such that

$$1001 HS_{t+1} \sim \mathcal{N}(\hat{\mu}, \hat{\Sigma}) \quad (28)$$

1002 where we use  $H = [I, I, \dots, I]^\top \in \mathbb{R}^{n_S \cdot E \times n_S}$  to project  $S_{t+1}$  to the dimension of the joint  $\hat{S}$ .

1003 We define the maximum likelihood loss

$$1004 \mathcal{L}(S_{t+1}) = p(\hat{S} | S_{t+1}) = \frac{1}{|2\pi\hat{\Sigma}|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\hat{S} - HS_{t+1}) \hat{\Sigma}^{-1} (\hat{S} - HS_{t+1}) \right) \quad (29)$$

1005 such that

$$1006 \log(\mathcal{L}(S_{t+1})) = -\frac{1}{2} \log(|2\pi\hat{\Sigma}|) - \frac{1}{2} (\hat{S} - HS_{t+1}) \hat{\Sigma}^{-1} (\hat{S} - HS_{t+1}). \quad (30)$$

1007 We aim to obtain the maximizer of the log-likelihood such that

$$1008 \bar{S}_{t+1} = \arg \max_{S_{t+1}} \log(\mathcal{L}(S_{t+1})). \quad (31)$$

1009 Consequently,

$$1010 \begin{aligned} \frac{\partial}{\partial S_{t+1}} \log(\mathcal{L}(S_{t+1})) &= -\frac{1}{2} H^\top \hat{\Sigma}^{-1} (\hat{S} - HS_{t+1}) := 0 \\ &\Rightarrow H^\top \hat{\Sigma}^{-1} \hat{S} - H^\top \hat{\Sigma}^{-1} H \bar{S}_{t+1} = 0 \\ &\Rightarrow \bar{S}_{t+1} = \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{S}. \end{aligned} \quad (32)$$

1011 As a result, we obtain

$$1012 \bar{\mu} = \mathbb{E}[\bar{S}_{t+1}] = \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{\mu} \quad (33)$$

and

$$\begin{aligned}
\bar{\Sigma} &= \text{Var} [\bar{S}_{t+1}] = \text{Var} \left[ \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{S} \right] \\
&= \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \text{Var} [\hat{S}] \hat{\Sigma}^{-1} H \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} \\
&= \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} H^\top \hat{\Sigma}^{-1} \hat{\Sigma} \hat{\Sigma}^{-1} H \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1} = \left( H^\top \hat{\Sigma}^{-1} H \right)^{-1}
\end{aligned} \tag{34}$$

which corresponds to standard results in Kalman fusion.

To estimate  $\Sigma^\Delta$ , we interpret  $\{\hat{\mu}^e\}_{e=1}^E$  as samples from a distribution whose mean is known to be  $\bar{\mu}$ . With this, the maximum likelihood estimate of  $\Sigma^\Delta$  can be obtained as

$$\bar{\Sigma}^\Delta = \frac{1}{E} \sum_{e=1}^E (\hat{\mu}^e - \bar{\mu}) (\hat{\mu}^e - \bar{\mu})^\top.$$

However, as the cross-correlations  $\hat{\Sigma}^{ef}$  are unknown in practice, we approximate the Kalman fusion results (33) and (34) using covariance intersection fusion Julier & Uhlmann (2001) with uniform weights, making use of Assumption 1. This results in

$$\bar{\Sigma} = \left( \frac{1}{E} \sum_{e=1}^E \left( \hat{\Sigma}^e \right)^{-1} \right)^{-1} \tag{35}$$

and

$$\bar{\mu} = \bar{\Sigma} \left( \frac{1}{E} \sum_{e=1}^E \left( \hat{\Sigma}^e \right)^{-1} \hat{\mu}^e \right). \tag{36}$$

### D.3 INFORMATION LOSS ALONG A ROLLOUT

As introduced in (4.3), the total information loss incurred during a equals the accumulated entropy of the state: Hence, we can estimate the environment state as

$$\underline{1,2,\dots,T} \underline{S_0} \bar{S}_{t+1} = \underline{s_0} \bar{\mu}(S_t, \underline{A_0} = \underline{a_0} A_t) + \underline{\bar{L}}(S_t, A_t) W_t, \underline{1=1 \dots T} = \sum_{t=0}^{T-1} \tag{37}$$

with  $\underline{\bar{L}} \underline{\bar{L}}^\top = \bar{\Sigma}$  and  $W_t \sim \mathcal{N}(0, I)$ . □

#### D.2.1 PROOF OF LEMMA 4.3

*Proof.*

$$\begin{aligned}
&\mathbb{H} \left( \bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T \right) \\
&= \sum_{t=0}^{T-1} \mathbb{H} \left( \bar{S}_{t+1} | \bar{S}_1, \bar{S}_2, \dots, \bar{S}_t, S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T \right) \\
&\stackrel{(a)}{=} \sum_{t=0}^{T-1} \mathbb{H} \left( \bar{S}_{t+1} | \bar{S}_1, \bar{S}_2 \dots \bar{S}_t, S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T \right) \\
&\stackrel{(b)}{=} \sum_{t=0}^{T-1} \mathbb{H} \left( \bar{S}_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1} \right) \\
&= \sum_{t=0}^{T-1} \mathbb{H} \left( \tilde{S}_{t+1} \right)
\end{aligned}$$

We continue here using the quantities we estimated in the previous section. To estimate  $\Sigma^\Delta$ , we interpret  $\{\hat{\mu}^1\}_{e=1}^E$  as samples from a distribution whose mean is known to be  $\bar{\mu}$ . With this, the maximum likelihood estimate of  $\Sigma^\Delta$  can be obtained trivially as

$$\bar{\Sigma}^\Delta = \frac{1}{E} \sum_{e=1}^E (\hat{\mu}^e - \bar{\mu}) (\hat{\mu}^e - \bar{\mu})^\top. \quad (38)$$

where (a) follows from causality and (b) follows from the Markov property.  $\square$

### D.3 INFOPROP STATE

As introduced in (5)(15), the Infoprop state is described as defined as

$$\tilde{S}_{t+1} := \mathbb{E} \left[ \bar{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1} \right] = \mathbb{E}_{\tilde{P}_{S,IP}} \left[ S_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}, U_t \right] \quad (39)$$

with mean  $\tilde{\mu}(s_t, a_t, \hat{s}_{t+1}) = \bar{\mu}(s_t, a_t) + K(s_t, a_t) (\bar{L}(s_t, a_t) w_t + \bar{L}^\Delta(s_t, a_t) n_t)$ , variance  $\tilde{\Sigma}(s_t, a_t, \hat{s}_{t+1}) = (I - K(s_t, a_t)) \bar{\Sigma}(s_t, a_t)$ ,  $K(s_t, a_t) = \bar{\Sigma}(s_t, a_t) (\bar{\Sigma}(s_t, a_t) + \bar{\Sigma}^\Delta(s_t, a_t))^{-1}$ ,  $\tilde{L}(s_t, a_t) \tilde{L}(s_t, a_t)^\top = \tilde{\Sigma}(s_t, a_t)$ , and  $\tilde{L}^\Delta(s_t, a_t) \tilde{L}^\Delta(s_t, a_t)^\top = \bar{\Sigma}^\Delta(s_t, a_t)$ .

Combining (3)(11) and Assumption 2, we have

$$\hat{S}_{t+1} = \check{S}_{t+1} + L^\Delta(S_t, A_t) N_t. \quad (40)$$

and according to (3)

$$\check{S}_{t+1} = \mu(s_t, a_t) + L(s_t, a_t) W_t.$$

Plugging the respective maximum likelihood estimates into (40) yields

$$\hat{S}_{t+1} = \bar{S}_{t+1} + \bar{L}^\Delta(S_t, A_t) N_t \quad (41)$$

We with

$$\bar{S}_{t+1} = \bar{\mu}(S_t, A_t) + \bar{L}(S_t, A_t) W_t \quad (42)$$

according to (13). As we can generally consider model uncertainty as independent from process noise, i.e.  $N_t \perp W_t$ , such that the Infoprop state

$$\tilde{S}_{t+1} = \mathbb{E} \left[ \check{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1} \right] = \mathbb{E}_{\tilde{P}_{S,IP}} \left[ S_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}, U_t \right]$$

$\tilde{S}_{t+1} = \mathbb{E}[\check{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}]$  can be computed using a standard Kalman update.

The general form of the Kalman update Simon (2006) considers two Gaussian RVs  $X \sim \mathcal{N}(\mu_X, \Sigma_X)$  and  $Y = X + N$  with  $N \sim \mathcal{N}(0, \Sigma_N)$  and  $X \perp N$ . Then, given an observation  $y$  we can compute the conditional expectation of  $X$

$$\mathbb{E}[X | Y = y] \sim \mathcal{N}(\mu_{X|Y=y}, \Sigma_{X|Y=y}) \quad (43)$$

with  $\mu_{X|Y=y} = \mu_X + K(y - \mu_X)$  and  $\Sigma_{X|Y=y} = (I - K) \Sigma_X$ , where  $K = \Sigma_X (\Sigma_X + \Sigma_N)^{-1}$ . Following (D.3),  $X$  and  $N$  represent the maximum likelihood estimates of the environment state

$$\mu_{X|Y=y} = \mu_X + K(y - \mu_X), \quad (44)$$

$$\Sigma_{X|Y=y} = (I - K) \Sigma_X, \quad (45)$$

and

$$K = \Sigma_X (\Sigma_X + \Sigma_N)^{-1}. \quad (46)$$

Following (15), we can compute the Infoprop state via (43) choosing

$$\mu_X = \bar{\mu}(s_t, a_t), \quad (47)$$

$$\Sigma_X = \bar{\Sigma}(s_t, a_t), \quad (48)$$

$$\Sigma_N = \bar{\Sigma}^\Delta(s_t, a_t), \quad (49)$$

and

$$y = \bar{\mu}(s_t, a_t) + \bar{L}(s_t, a_t)w_t + \bar{L}^\Delta(s_t, a_t)n_t. \quad (50)$$

This yields the propagation equation of the Infoprop state

$$\tilde{S}_{t+1} = \tilde{\mu}(S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}) + \tilde{L}(S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1})U_t \quad (51)$$

with

$$\tilde{\mu}(s_t, a_t, \hat{s}_{t+1}) = \bar{\mu}(s_t, a_t) + K(s_t, a_t) (\bar{L}(s_t, a_t)w_t + \bar{L}^\Delta(s_t, a_t)n_t), \quad (52)$$

$$\tilde{\Sigma}(s_t, a_t, \hat{s}_{t+1}) = (I - K(s_t, a_t)) \bar{\Sigma}(s_t, a_t), \quad (53)$$

$$K(s_t, a_t) = \bar{\Sigma}(s_t, a_t) (\bar{\Sigma}(s_t, a_t) + \bar{\Sigma}^\Delta(s_t, a_t))^{-1}, \quad (54)$$

$$\tilde{L}(s_t, a_t) \tilde{L}(s_t, a_t)^\top = \tilde{\Sigma}(s_t, a_t), \quad (55)$$

and

$$\bar{L}^\Delta(s_t, a_t) \bar{L}^\Delta(s_t, a_t)^\top = \bar{\Sigma}^\Delta(s_t, a_t). \quad (56)$$

#### D.4 INDUCED STATE DISTRIBUTION BY THE INFOPROP ROLLOUT

**Lemma 3.** *As introduced in (57), the next state distribution induced by the Infoprop rollout is the same as that given by the estimated ground truth:*

$$\tilde{S}_{t+1} \stackrel{\text{dist}}{=} \bar{S}_{t+1} \quad (57)$$

*Proof.* We show equality in distribution via comparison of the cumulative distribution functions (CDF) of  $\tilde{S}_{t+1}$  and the epistemic noise  $\bar{S}_{t+1}$ . If we can show that the CDFs are identical, i.e.  $\mu_X = \bar{\mu}(s_t, a_t)$ ,  $\Sigma_X = \bar{\Sigma}(s_t, a_t)$ , and  $\Sigma_N = \bar{\Sigma}^\Delta(s_t, a_t)$ . Further,  $y$  represents a sample from the model  $y = \bar{\mu}(s_t, a_t) + \bar{L}(s_t, a_t)w_t + \bar{L}^\Delta(s_t, a_t)n_t$ . Plugging these into (43), yields (5).  $\mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) \quad \forall \bar{s}_{t+1} \in \mathcal{S}$ , the equality in distribution follows.

We compute  $\mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1})$  using  $\tilde{S}_{t+1} = \mathbb{E}[\tilde{S}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}]$  and marginalizing over  $\hat{S}_{t+1}$ .

$$\mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \int_{\mathcal{S}} \mathbb{P}(\mathbb{E}[\tilde{S}_{t+1} | \hat{S}_{t+1}] \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) f_{\hat{S}_{t+1}}(\hat{s}_{t+1}) d\hat{s}_{t+1} \quad (58)$$

with  $f_{\hat{S}_{t+1}}$  the probability density function of  $\hat{S}_{t+1}$ .

By construction,  $\mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1}]$  describes the behavior of  $\bar{S}_{t+1}$  given  $\hat{S}_{t+1}$ . Consequently,

$$\mathbb{P}(\mathbb{E}[\bar{S}_{t+1} | \hat{S}_{t+1}] \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) = \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) \quad (59)$$

and therefore

$$\mathbb{P}(\tilde{S}_{t+1} \leq \bar{s}_{t+1}) = \int_{\mathcal{S}} \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) f_{\hat{S}_{t+1}}(\hat{s}_{t+1}) d\hat{s}_{t+1}. \quad (60)$$



The right hand side of (60) represents the law of total probability for  $P(\bar{S}_{t+1} \leq \bar{s}_{t+1})$ .

$$\mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) = \int_{\mathcal{S}} \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1} | \hat{S}_{t+1} = \hat{s}_{t+1}) f_{\hat{S}_{t+1}}(\hat{s}_{t+1}) d\hat{s}_{t+1}. \quad (61)$$

Therefore, we have

$$\mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) = \mathbb{P}(\bar{S}_{t+1} \leq \bar{s}_{t+1}) \quad \forall \bar{s}_{t+1} \in \mathcal{S} \quad (62)$$

and can conclude

$$\bar{S}_{t+1} \stackrel{\text{dist}}{=} \bar{S}_{t+1}. \quad (63)$$

□

#### D.5 INFORMATION LOSS ALONG A INFOPROP ROLLOUT

**Lemma 4.** *As introduced in (4.3), the total information loss incurred during a Infoprop equals the accumulated entropy of the Infoprop state:*

$$\mathbb{H}(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T) = \sum_{t=0}^{T-1} \mathbb{H}(\tilde{S}_{t+1}) \quad (64)$$

*Proof.*

$$\begin{aligned} & \mathbb{H}(\bar{S}_1, \bar{S}_2, \dots, \bar{S}_T | S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T) \\ &= \sum_{t=0}^{T-1} \mathbb{H}(\bar{S}_{t+1} | \bar{S}_1, \bar{S}_2, \dots, \bar{S}_t, S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1 \dots \hat{S}_T = \hat{s}_T) \\ &\stackrel{(a)}{=} \sum_{t=0}^{T-1} \mathbb{H}(\bar{S}_{t+1} | \bar{S}_1, \bar{S}_2 \dots \bar{S}_t, S_0 = s_0, A_0 = a_0, \hat{S}_1 = \hat{s}_1, \dots, \hat{S}_T = \hat{s}_T) \\ &\stackrel{(b)}{=} \sum_{t=0}^{T-1} \mathbb{H}(\bar{S}_{t+1} | S_t = s_t, A_t = a_t, \hat{S}_{t+1} = \hat{s}_{t+1}) \\ &= \sum_{t=0}^{T-1} \mathbb{H}(\tilde{S}_{t+1}) \end{aligned} \quad (65)$$

where (a) follows from causality and (b) follows from the Markov property.

□

## E EXPERIMENTS

### E.1 EXPERIMENTAL SETUP

We used Weights&Biases<sup>4</sup> for logging our experiments and run 5 random seeds per experiment.

The respective hyperparameters for Infoprop-Dyna on MuJoCo are given below. Table 2 addresses model learning, Table 3 the Infoprop mechanism, and Table 4 training the model-free agent.

Table 2: Hyperparameters used to train the model of Infoprop-Dyna in the Mujoco Tasks.

| Hyperparameter              | Halfcheetah           | Walker | Hopper | Ant     |
|-----------------------------|-----------------------|--------|--------|---------|
| ensemble size $E$           | 7                     |        |        |         |
| number of hidden neurons    | 200                   |        |        | 400     |
| number of hidden layers     | 4                     |        |        |         |
| learning rate               | 0.0003                | 0.0006 | 0.0004 | 0.001   |
| weight decay                | 0.00005               | 0.0007 | 0.0008 | 0.00002 |
| patience for early-stopping | 10                    | 9      | 8      | 9       |
| retrain interval            | 250 environment steps |        |        |         |

<sup>4</sup><https://wandb.ai/site>

Table 3: Hyperparameters of the Infoprop rollouts in the Mujoco Tasks.

| Hyperparameter                            | Halfcheetah           | Walker | Hopper | Ant |
|---|-----------------------|--------|--------|-----|
| accurate quantile $\zeta_1$               | 0.99                  |        |        |     |
| exceptionally accurate quantile $\zeta_2$ | 0.01                  |        |        |     |
| scaling factor $\xi$                      | 100                   |        |        |     |
| rollout interval                          | 250 environment steps |        |        |     |
| rollout batch size                        | 100000                |        |        |     |

Table 4: Hyperparameters used to train the SAC agent of Infoprop-Dyna in the Mujoco Tasks.

| Hyperparameter           | Halfcheetah | Walker | Hopper | Ant    |
|--------------------------|-------------|--------|--------|--------|
| number of hidden neurons | 1024        |        | 512    | 1024   |
| number of hidden layers  | 2           |        |        |        |
| learning rate            | 0.0003      | 0.0002 | 0.0004 | 0.0005 |
| SAC target entropy       | -6          | -7     | 1      | 0      |
| target update interval   | 1           | 4      | 6      | 5      |
| update steps $G$         | 10          |        |        | 20     |

The results for SAC, MBPO and MACURA are obtained from Frauenknecht et al. (2024).

## E.2 PREDICTION QUALITY

We provide additional results for the rollout consistency experiments introduced in Section 6.2. Figure 6 depicts model-based rollouts for the 10<sup>th</sup> dimension of hopper under MBPO, MACURA and Infoprop-Dyna when setting the maximum rollout length of all approaches to 100. In the original experiment depicted in Figure 4b the maximum rollout length was 11 for MBPO and 10 for MACURA, following the hyperparameter settings reported in the respective publications Janner et al. (2019); Frauenknecht et al. (2024).

We observe a vastly spread distribution of MBPO rollouts, as every rollout is propagated for 100 steps, irrespective of model uncertainty, as long as it does not reach a terminal state of the hopper task. MACURA rollouts have an improved consistency compared to MBPO, especially in the beginning of the rollouts. Over long horizons, however, the TS propagation mechanism and the single-step termination criterion cannot produce consistent data. In contrast, Infoprop-Dyna is able to propagate consistent rollouts over long horizons.

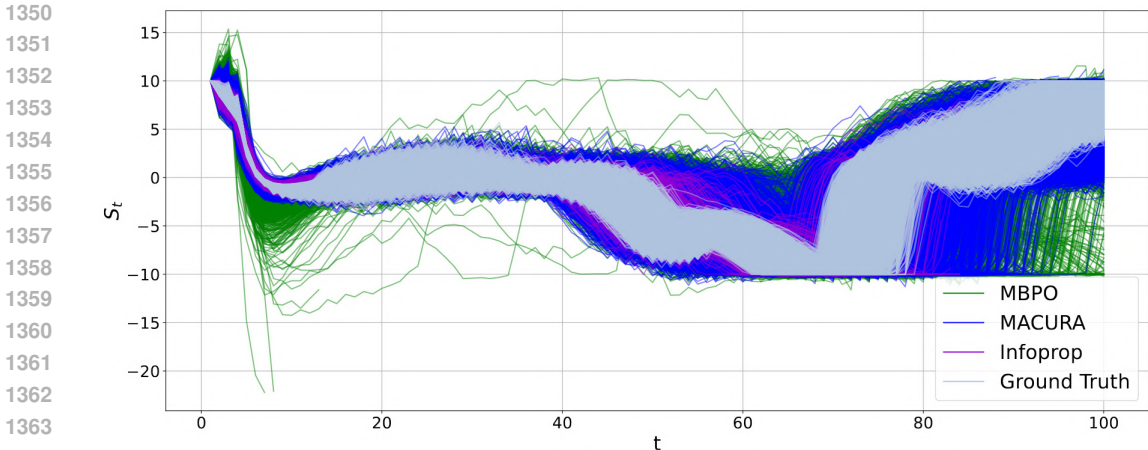


Figure 6: Rollout consistency MBPO vs. MACURA vs. Infoprop-Dyna for 100 steps. Comparison of the respective rollout mechanisms similar to Figure 4b but with a maximum rollout length of 100 for all approaches.

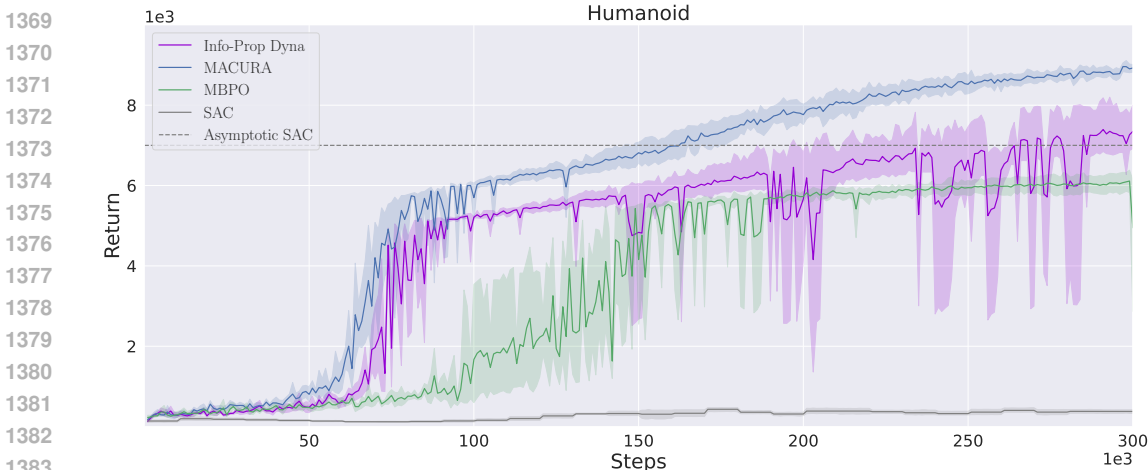


Figure 7: Performance on Humanoid

E.3 LIMITATIONS PERFORMANCE ON HUMANOID

Figure 7 depicts the return on MuJoCo humanoid. We observe instabilities in the performance of Infoprop-Dyna towards the end of training. We assume this occurs due to overfitting and plasticity loss in the critic of the model-free learner Nikishin et al. (2022); D’Oro et al. (2023). This is reflected in the peaking critic loss depicted in Figure 8 concurrently with the performance drops. We set the update ratio  $G$  (see Algorithm 3) to a relatively low value of 10 which explains the slower learning behavior than MACURA. For higher values of  $G$ , instabilities occur even earlier in the training process, underscoring our assumption of overfitting critics. A similar observation is reported in Frauenknecht et al. (2024), where low values of the scaling factor  $\xi$ , corresponding to accurate model rollouts, led to instabilities in learning.

Model rollout inconsistency does not appear to be the destabilizing factor, as rollout data is consistent with the environment distribution as depicted in Figure 10 and the rollout adaption mechanism seems to react to policy shifts induced by high critic losses through reducing the average rollout length as depicted in Figure 10.

E.4 INVESTIGATING INSTABILITIES IN LEARNING

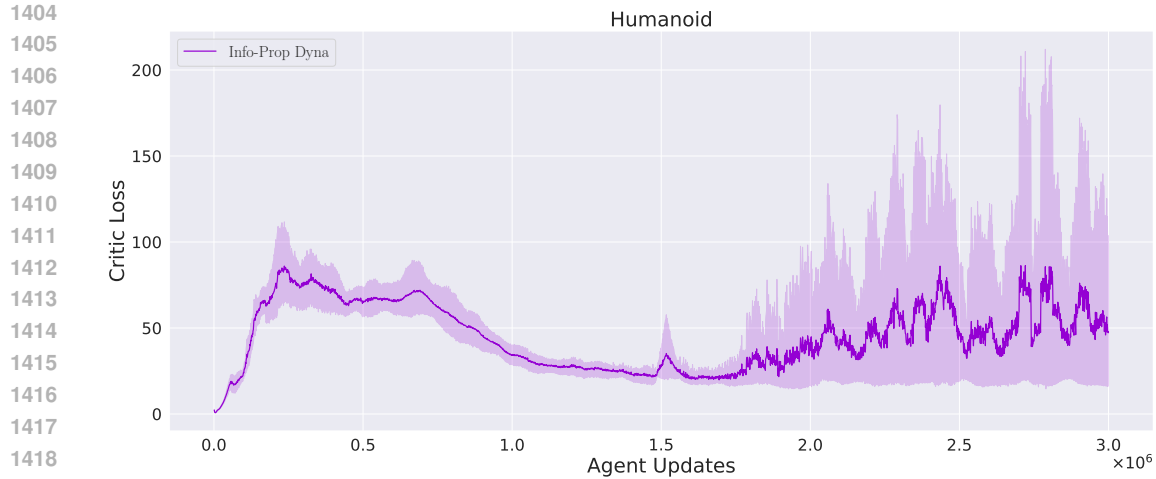
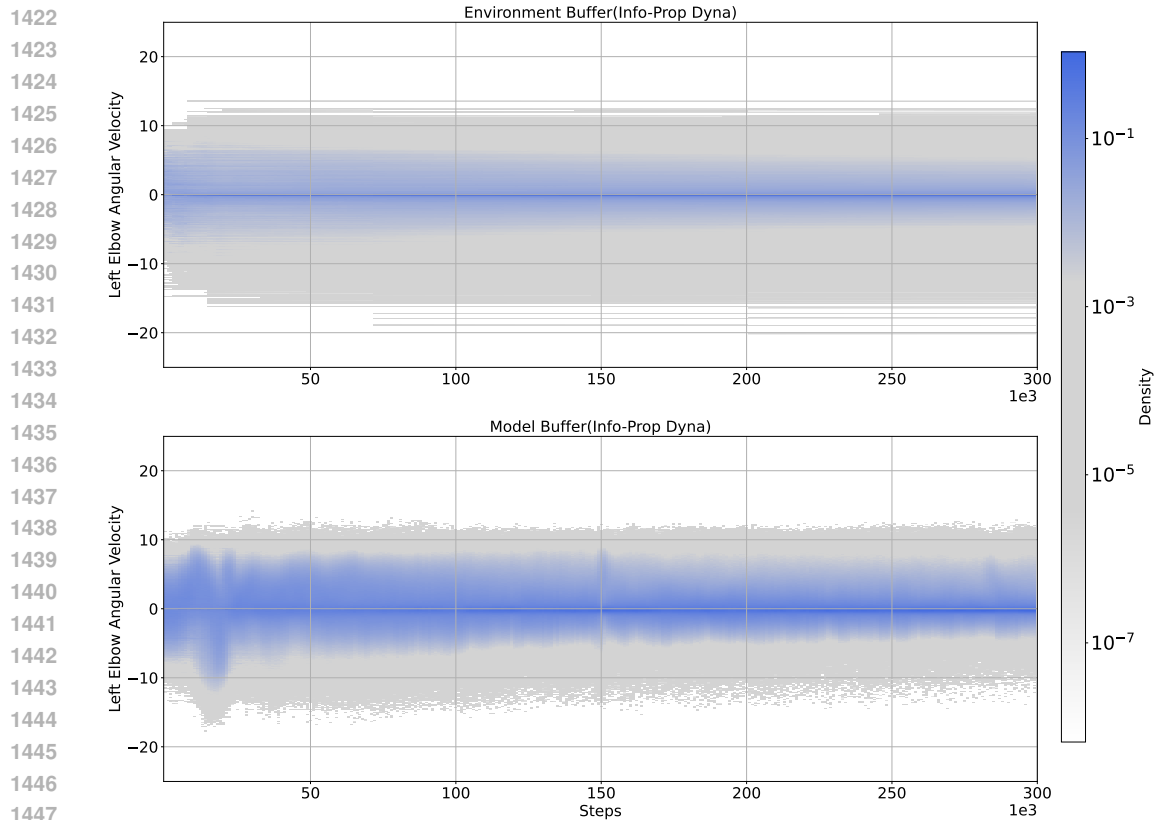


Figure 8: Critic Loss on Humanoid

Figure 9: Comparison between  $\mathcal{D}_{env}$  and  $\mathcal{D}_{mod}$  for the 45<sup>th</sup> dimension of Humanoid

1451

1452

1453

1454

1455

1456

1457

Although Infoprop gives better quality data over longer rollout horizons than TS rollouts, we observe instabilities in learning when naively integrating Infoprop into the conventional Dyna setting. We hypothesize that the main cause of these instabilities is due to the agent overfitting to the higher quality data produced by Infoprop rollouts, followed by loss of plasticity Nikishin et al. (2022); D’Oro et al. (2023). To investigate this, we carried out an ablation by varying the values of  $\zeta_1$ , which we introduced in Equation 22. This hyperparameter controls the size of the subset  $\mathcal{E}$  where the model is considered sufficiently accurate. The smaller the value of  $\zeta_1$ , the more aggressive the filtering of single-step information losses, leading to a smaller  $\mathcal{E}$ .

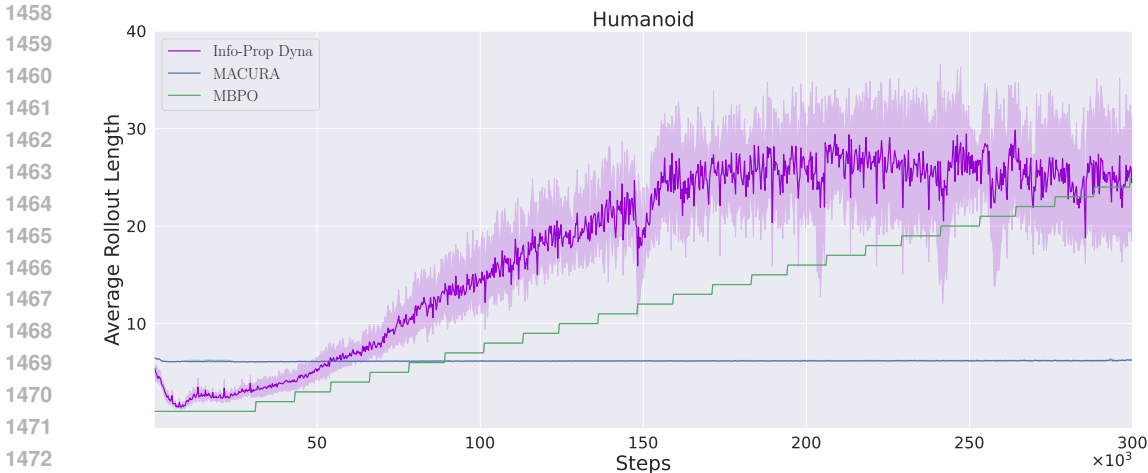


Figure 10: Average Rollout Length on Humanoid

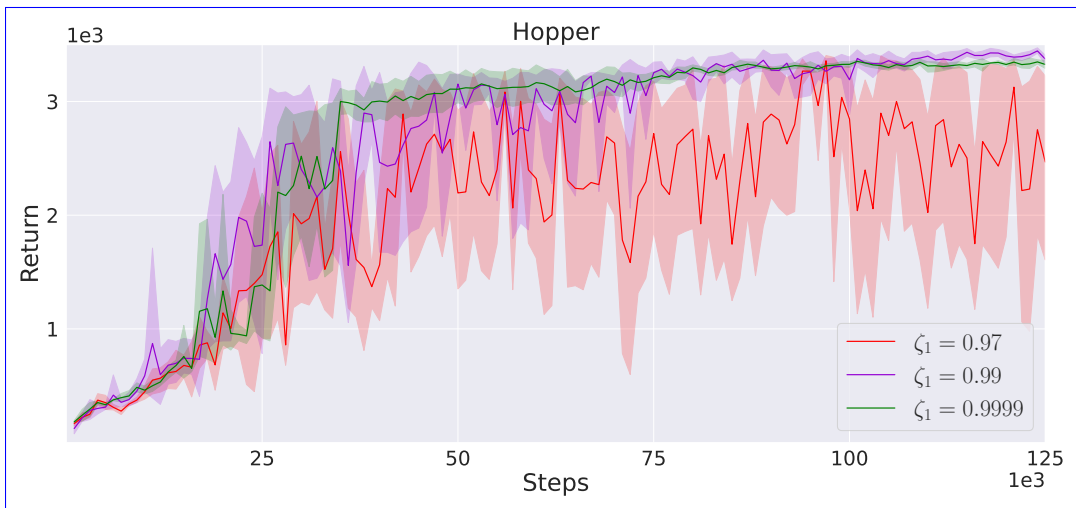


Figure 11: Ablation Study on Hopper.

Figure 11 shows the returns obtained on the Hopper task for three values of  $\zeta_1$ . For  $\zeta_1 = 0.97$ , we see that the returns are unstable throughout training, even though this setting gives the best quality data. On the other hand,  $\zeta_1 = 0.9999$  produces a more stable learning curve compared to  $\zeta_1 = 0.99$ , which was used for all the experiments in Section 6. This shows that better data quality does not necessarily lead to better training performance since if that was the case,  $\zeta_1 = 0.97$  would have produced the best performance. A similar observation is reported in Frauenknecht et al. (2024), where low values of the scaling factor  $\xi$ , corresponding to accurate model rollouts, led to instabilities in learning.

Our observations show that producing high-quality synthetic data in the conventional Dyna setting leads to issues seen in MFRL when using a high update-to-data (UTD) ratio. There have been recent works on regularization methods to counteract agent overfitting and loss of plasticity. One such approach is applying layer normalization Smith et al. (2023); Nauman et al. (2024). Figure 12 shows the same settings as in Figure 11 but with layer normalization applied to the critic and actor networks. It can be seen that even for  $\zeta_1 = 0.97$ , the learning is stable.

The primary aim of this paper is to introduce the conceptual framework of the Infoprop rollout, as well as show its application to MBRL. Hence, we do not spend additional effort on tuning the hyperparameters or adding regularizations since this takes us away from the main objective. We defer such enhancements for future work.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

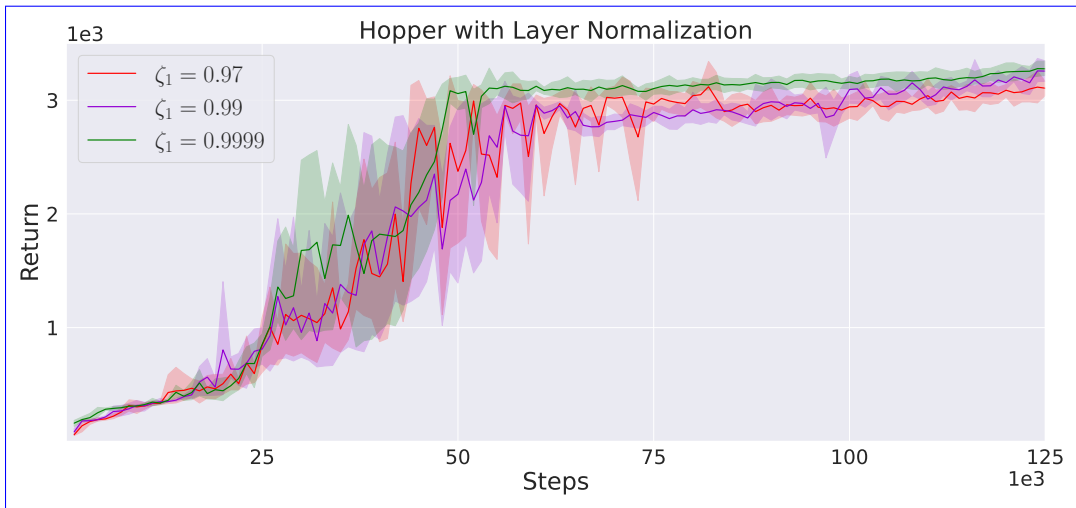


Figure 12: [Ablation Study on Hopper with layer normalization.](#)