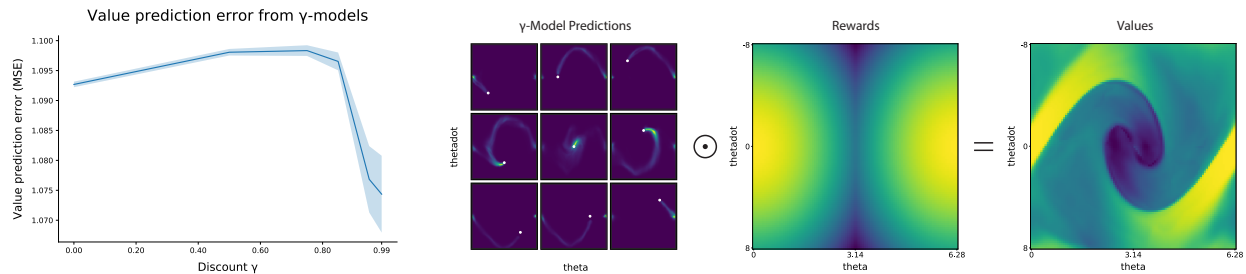


1 We thank the reviewers for their insights and suggestions. Questions primarily concerned the utility of  $\gamma$ -model  
 2 predictions and associated tradeoffs, policy-dependent models, and scalability. We clarify how  $\gamma$ -model predictions  
 3 are used for decision-making (with further elaboration on value estimation), discuss how the tradeoff induced by the  
 4 choice of  $\gamma$  is one that must be made in any RL algorithm, and comment on scaling to higher-dimensional environments.  
 5 Answers below will be included in expanded discussions in future versions of the paper.

6 **(R3) Utility of “jumpy” predictions.** One of the insights of this paper is that model-based RL algorithms do not  
 7 need to know the time at which a state will be encountered as long as states are sampled according to the discounted  
 8 occupancy. This does mean that decisions should be based on expectations over many samples (Figure R1b) and not  
 9 based on a single sample. In the case of R3’s car example, as long as states from 10 steps into the future are sampled  
 10 according to a probability of  $p(t = 10) = (1 - \gamma)\gamma^9$ , crashing at  $t = 10$  will be reflected in such an expectation.

11 **(R2, R3) Clarifications on RL experiments.** (i) In RL experiments,  $\gamma$ -models were **not** pre-trained; they were trained  
 12 online. (ii) Though bounds exist, it is difficult to relate control performance to state prediction accuracy in practice  
 13 [Lambert et al., 2020]. The intermediate quantity of value prediction error (Figure R1a) is more informative in this  
 14 respect, and shows a small accuracy improvement for  $\gamma > 0.95$ , consistent with the performance of  $\gamma$ -MVE vs MVE.



**Figure R1:** (a) Value prediction error for an occupancy with target discount  $\tilde{\gamma} = 0.99$  as a function of model discount  $\gamma$ . Higher  $\gamma$ 's perform slightly better. (b) Visual depiction of value estimation (in Pendulum) as an expectation of reward over  $\gamma$ -model predictions.

15 **(R1, R2) Tradeoff with  $\gamma$ .** The  $\gamma$ -model bridges the gap between model-based and model-free learning, and as such  
 16 allows one to trade off between training-time model-free errors (referred to as bootstrap error accumulation; Kumar et al.  
 17 2019) and testing-time model-based rollout errors. The tradeoff between these two is not a weakness of the  $\gamma$ -model,  
 18 nor is it even unique to the  $\gamma$ -model; such compounding errors are unavoidable in sequential decision-making. The  
 19 advantage of the  $\gamma$ -model is that it allows for more careful interpolation between these extremes.

20 **(R2, R3) Policy-dependence.** The  $\gamma$ -model is policy-dependent in the same way that a Q-function is. While the single-  
 21 step transition distribution is certainly independent of a policy, parametric single-step models are still policy-dependent  
 22 because the policy determines the training data distribution for statistical learning.

23 **(R1, R4) Scaling to image-based environments.** We agree that experiments in more complex domains would improve  
 24 the submission. Unfortunately, model-based experiments in image domains like Atari games are quite complex to set  
 25 up, generally requiring specialized image prediction models that can take weeks to train [Kaiser et al., 2019]. We have  
 26 found that  $\gamma$ -model training scales to state spaces of the MuJoCo gym benchmark suite for discounts up to  $\gamma = 0.9$ , but  
 27 leave extensions to image-based environments for future work.

28 **Other questions from Reviewer 2** ① *Gradients:* The gradient computation indeed does not pass through the  
 29 bootstrapped target. This is discussed in L211-L215 in Section 6 “Practical Training of  $\gamma$ -Models”.

30 **Other questions from Reviewer 3** ① *Prediction visualization:* The first five columns of Figure 3 depict learned  
 31  $\gamma$ -model predictions. The only Monte Carlo trajectory estimates are in the final column for comparison. ② *Relation*  
 32 *to n-step models:* Training an  $n$ -step model for  $n$  larger than 10 or 20 quickly becomes impractical, and is impossible  
 33 without on-policy samples. By directly learning the discounted occupancy, we are able to train a model with much  
 34 larger effective (probabilistic) horizon while only using off-policy single-step transitions. ③ *Bellman equation:* There  
 35 is no recursion in Equation 2 because it describes a bootstrapped target distribution (akin to a bootstrapped target value  
 36  $r(s_t) + \gamma V(s_{t+1})$ ) and is not a Bellman equation by itself. ④ *Required rollout lengths:* Figure 2b shows the length  
 37 of a  $\gamma$ -model-based rollout required to recover 95% of the probability mass of an occupancy with discount  $\tilde{\gamma}$ . This is  
 38 more precise than an effective horizon. For example, the effective horizon for  $\tilde{\gamma} = 0.99$  with a single-step model is 100,  
 39 which recovers only  $1 - .99^{101} \approx 64\%$  of the full infinite-horizon occupancy. ⑤ *Energy-based model:* We introduce  
 40 the EBM formulation because it provides the simplest way to explain how temporal difference updates may be used to  
 41 train generative models. A neural generator reduces training times, but this is more an implementation choice than a  
 42 fundamental component of infinite-horizon prediction, so we opt to first explain the framework with as few moving  
 43 pieces as possible. Similarly, we have found that the  $\gamma$ -model may also be instantiated as a normalizing flow, and will  
 44 add experiments comparing the flow and GAN instantiations of  $\gamma$ -models.

45 **References** N Lambert et al, *arXiv:2002.04523*. A Kumar et al, *arxiv:1906.00949*. L Kaiser et al, *arXiv:1903.00374*.