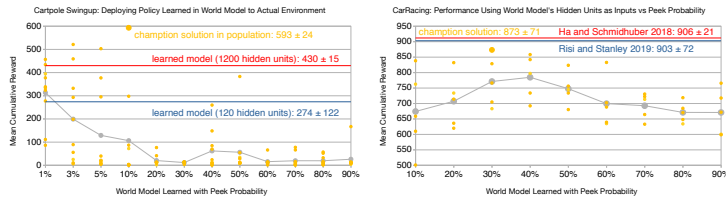


1 We thank the reviewers for their feedback and comments on the manuscript! Indeed, this work was largely an exploration
2 of the idea: "Can a world model be learned without using a forward-predictive loss?" Owing to the intense research
3 effort spent on learning forward predictive models in the presence of many a flavor of forward-predictive auxiliary
4 losses, we focused our efforts on simple tasks that best illustrated the trade-offs of our exploration. We are assured to
5 see that we have achieved this as all 3 reviewers thought our method was interesting, and that our presentation was clear.
6 While adding more complex test environments would certainly strengthen our work, we assert that the experiments
7 presented here provide a compelling sketch for a different way to train world models.

8 Because of the relative maturity of the model-based literature, we were initially hesitant to make direct comparisons
9 with model-based baselines because doing so immediately invites a comparison that we did not intend to draw. However,
10 we fully agree that including such a baseline is useful, informative, and, if nothing else, establishes a long term goal for
11 this program. We have now included model-based baselines for the swing-up cartpole and car-racing tasks (see below),
12 and we have launched more detailed baseline experiments for car-racing (suggested by R3) and gridworld.



Cartpole Swingup	Reward
Explicitly learned model from data (120 hidden units)	274 ± 122
Explicitly learned model from data (1200 hidden units)	430 ± 15
Champion solution in our population (120 hidden units)	593 ± 24

CarRacing-v0	Reward
Ha and Schmidhuber (NeurIPS 2018)	906 ± 21
Risi and Stanley (GECCO 2019)	903 ± 72
Champion solution in our population	873 ± 71

13 The full baselines and discussion will be readily available in the revision. To our surprise, we find it interesting that our
14 approach can produce models that outperform an explicitly learned model with the same architecture size (120 units) for
15 cartpole transfer task, but this advantage goes away if we scale up the learned model size by 10x. For car-racing, only
16 the best solution out of a population of trials gets near the performance of population-based baselines from the literature,
17 but honestly we did not expect our new implicit approach of learning world models to be this close in performance.

18 With respect to R1, R2 and R3's comments on prior work (especially the Predictron), we of course will add these
19 overlooked references, and provide a discussion of similarities and differences of these works. In particular for the
20 Predictron, both methods are after the same goals: a latent model emerging from optimization. The predictron does
21 this via learning representations, dynamics, and value estimates end-to-end over imagined rollouts of varying lengths.
22 The primary innovation of observational dropout is to recouple the "imagined" rollouts back to the real observation
23 space—i.e., any sequence of transitions within our learned world models always end back on a "real" frame from
24 the environment, with "real" environment steps in the intervening frames. Additionally, where the predictron derives
25 consistency by demanding that its value estimates obey a Bellman equation, we only demand that the learned dynamics
26 model facilitate the learning of a policy. While similar in spirit, we hope the difference here is clear—it's interesting,
27 for example, to consider how a predictron could be combined with observational dropout, where at each step, instead of
28 necessarily seeing another ground-truth frame from the environment, the predictron sees its own 1-step prediction.

29 Regarding R1 and R2's concerns on generalization (whether to "harder" or "unseen" tasks): we agree that understanding
30 generalization is crucial, and regret that we did not have more space to discuss this point in the manuscript. In sum: for
31 high observational dropout, too much information is lost, and the world model does not have the capability to accurately
32 reconstruct the system due to noise. At low dropout, the world model might as well not be present, because the policy
33 can ignore the occasionally noisy frame. The extent to which the world model can then transfer to "unseen tasks" is
34 dependent on how well the world model fully captures the dynamics of the system. In some sense, world models learned
35 in this way are "bad" at generalizing, because they are trained to only learn the "relevant" dynamics of the task being
36 studied (as in the grid world task, where it only learned transitions in some but not all directions). In another sense,
37 they're "good" at generalizing, because the "relevant" dynamics of the current task might also be the relevant dynamics
38 of "unseen" tasks (i.e., as in balance cartpole, initialized at an angle not seen during training). It is difficult to make a
39 precise statement here, because the quality of the learned world model truly depends on the dynamics being studied.

40 R3 had an interesting comment regarding the relative power of the policy versus the world model—indeed, a sufficiently
41 powerful policy could extract meaningful information from the world model, even if the world model was not actually
42 forward predictive at all. In this limit, the joint policy + world model system starts looking like a strange sort of
43 recurrent network. We agree with this assessment, and will add a short discussion highlighting this limit. An alternative
44 interpretation is that our procedure is a way of encouraging a recurrent model's internal state to be dynamically
45 relevant/interpretable—first by mediating a policy's observation space via that lossy dynamics module, and then by
46 constraining the policy itself to be simple, so that it is not, in R3's words, "sufficiently powerful". We are absolutely
47 interested in pursuing this interpretation through an information bottleneck lens, but felt it was out of the scope of this
48 work. Finally, we absolutely will add more details on the number of training steps used for the gridworld tasks (an
49 equal number was used for both architectures, to answer R3), as well as details on the observation space of cartpole. We
50 hope this response has addressed the reviewer's concerns, and would appreciate a reevaluation of your scores!