

1 We thank the reviewers for their detailed comments and positive evaluation of our work. The questions primarily
 2 concerned (1) evaluation in an environment with a larger action space, (2) performance of ME-TRPO, (3) analysis of
 3 the gap between MBPO and the baselines, and (4) discussion of the tightness of the theoretical bound.

4 **(R1) Code release.** Code for reproducing our experiments is now
 5 available on GitHub. To preserve anonymity, we do not link directly to
 6 the repository.

7 **(R1) Larger environment.** We provide results on the Humanoid en-
 8 vironment, requested by R1, in Figure 1. We will add these to the final
 9 version.

10 **(R1) Intuitions for the theory.** We will expand the discussion in Sec-
 11 tion 4 to present a better intuition for the practical implications of the
 12 theory. The theory suggests that: (1) short-horizon rollouts may be ben-
 13 efiticial in some settings; (2) incorporating the model can allow for larger
 14 policy changes while still achieving monotonic improvement, but only
 15 if the model generalizes well to changes in the policy – the worst-case
 16 generalization does not achieve this, but we empirically find that real
 17 models on MuJoCo benchmark tasks generalize substantially better.

18 **(R1, R2) Analysis of comparative performance.** R1 and R2 asked
 19 about the sources of improvement of our method over the baselines.
 20 Here we elaborate on our choice of ablations and how they address this
 21 question.

- 22 1. The 500-length rollout ablation in the paper’s Figure 3 is the suggested ME-SAC baseline, as it uses the
 23 ensemble to generate model rollouts with lengths on the order of the task horizon for consumption by SAC. We
 24 conclude from this result that truncated rollouts are a primary source of the performance difference between
 25 our method and **ME-TRPO**.
- 26 2. To better understand the difference between using model data directly for training and for improved target
 27 value estimates (as done in **MVE** and **STEVE**), we implemented the value expansion technique on top of SAC
 28 to control for the underlying model-free algorithm. This comparison is found in the paper’s Figure 3.
- 29 3. We do not include a separate ablation for **PETS** because the comparison between MBPO and PETS is already
 30 well-controlled. The model ensembles are the same in both methods, so the difference in performance is
 31 attributable to the different ways the model is used: planning by sampling from a fixed prior in PETS and
 32 policy optimization in MBPO.

33 **(R2) ME-TRPO baseline.** R2 raised concerns about the relatively poor performance of ME-TRPO on the full-length
 34 HalfCheetah. The ME-TRPO paper evaluates on modified tasks, with horizons of 100 or 200, making their reported
 35 results not representative of the standard benchmarks. Our results use the authors’ code and are representative of the
 36 actual performance of the method. An independent benchmarking of model-based RL algorithms, released after the
 37 NeurIPS deadline, reported the same results from ME-TRPO on the full-length environments [Wang et al., 2019].

38 **(R2) Elaboration on branched rollouts.** The branching rollouts use the marginal distribution from a previous policy
 39 as an initial state distribution for truncated model rollouts. In practice, this amounts to sampling a state $s \sim \mathcal{D}$
 40 from the environment replay buffer, rolling out under the model for at most k steps using the current policy, and using these
 41 model predictions for policy optimization.

42 **(R3) Tightness of the bound.** Our bound is tightest in MDPs in which a single differing
 43 action or transition leads two trajectories to permanently diverge, as in the binary tree MDP
 44 in Figure 2. A crucial step in proving Theorem 1 is that if two agents select differing actions
 45 with ϵ probability, then their state marginals diverge by ϵt in total variation (Lemma B2). In
 46 Figure 2, the amount of divergence is exactly $1 - (1 - \epsilon)^t$, which is close to ϵt when ϵ is small.
 47 We are not aware of a way to create a tighter bound while still handling this pathological
 48 MDP. Similar proof techniques are used to analyze the TRPO and CPI algorithms.

49 **References**

50 Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang,
 51 Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning.
 52 *arXiv preprint arXiv:1907.02057*, 2019.

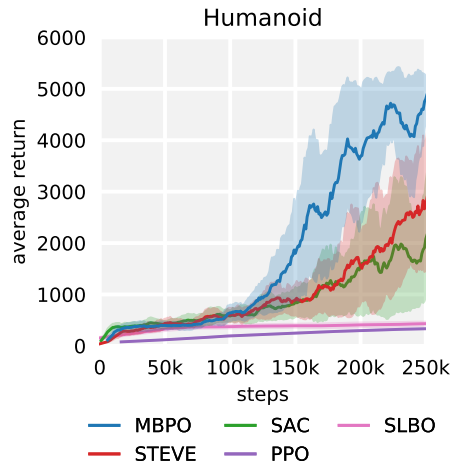


Figure 1: Results on Humanoid-v2. MBPO results are averaged over four seeds. The short rebuttal period did not allow for running all baselines to convergence, but we will add them to the final.

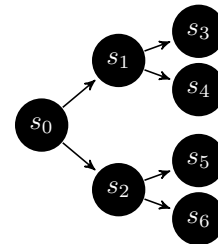


Figure 2: An example MDP where the bound is presented in Theorem 1 is nearly tight.