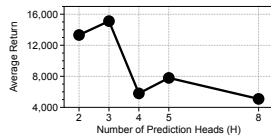


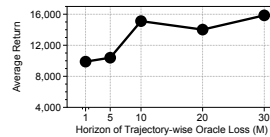
1 We thank all reviewers for carefully reading our paper and their valuable comments. We appreciate that our paper is
 2 recognized for several positive aspects: [R1, R2, R4] novel and well-motivated, [R2, R3, R4] clear write-up, and [All]
 3 extensive and strong experiments. Below are our responses to the reviewers, which we will incorporate in the final draft.



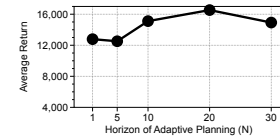
(a) Prediction heads

Mass	Head 1	Head 2	Head 3	Head 4	Head 5
0.25	0.0	0.0	0.0	100.0	0.0
0.50	0.1	0.0	0.0	0.0	99.9
1.50	56.3	43.7	0.0	0.0	0.0
2.50	0.0	44.3	55.7	0.0	0.0

(b) Trajectory assignment



(c) Trajectory-wise oracle loss



(d) Adaptive planning

4 [R2, R3, R4] **Effects of major hyperparameters.** We choose the number H of prediction heads and the horizon
 5 M of trajectory-wise oracle loss based on the trajectory assignments in training environments, i.e., how distinctively
 6 trajectories are assigned to prediction heads (see Figure 4 and 5 of the original draft and above Figure (b) for examples).
 7 Also, the horizon N of adaptive planning is set to the same value as M to match the horizon of trajectory-wise oracle
 8 loss and adaptive planning. We use the same hyperparameters H , M , and N for all environments. We will clarify the
 9 hyperparameter setups in the final draft.

10 Following the suggestions of R2, R3, and R4, we also perform ablation studies in HalfCheetah environments by varying
 11 the major hyperparameters of our method: $H \in \{2, 3, 4, 5, 8\}$, $M \in \{1, 5, 10, 20, 30\}$, and $N \in \{1, 5, 10, 20, 30\}$.
 12 Figure (a) shows that $H = 3$ achieves the best performance because three prediction heads are enough to capture
 13 the multi-modality of the training environments in our setting. When $H > 3$, the performance decreases because
 14 trajectories from similar dynamics are split into multiple heads as shown in Figure (b). However, as pointed out by R4,
 15 we expect that more heads would be effective in environments with more varying environmental factors. Figure (c) and
 16 (d) show that choosing the head per-trajectory is more effective and our method is robust to change in hyperparameters
 17 $M \geq 10$ and $N \geq 10$. We will include more comprehensive results for all environments in the final draft.

18 [R1] **Limitation of our method.** Thank you very much for your pointers. As we assume that MDPs with similar
 19 dynamics will behave similarly, the effectiveness of our method would be limited if the dynamics of unseen environments
 20 are significantly different from the dynamics of training environments. As you pointed out, gain from our method may
 21 decrease if the agent is not exposed to a variety of MDPs. We will clarify these limitations in the final draft.

22 [R1] **Clarification on problem formulation.** Our work addresses the dynamics generalization problem of model-
 23 based RL methods, where learned dynamics models fail to provide accurate predictions as the transition dynamics of
 24 environments change. Thank you for your suggestion, and we will clarify this in the final draft.

25 [R1] **Clarification on multi-modal nature of environments.** We assume that the transition dynamics distribution of
 26 MDP is multi-modal, which emerges as the environmental factors (e.g., mass and length of the agent) change. We
 27 remark that capturing this property is important as environmental factors are ever-changing in real-world environments.
 28 Prior model-based RL methods that do not consider this property fail to generalize as the future prediction from
 29 dynamics models becomes inaccurate. We will clarify this in the final draft.

30 [R2] **Editorial comment.** We will remove the “Due to space limitation” phrase in the final draft.

31 [R3] **Novelty.** As R1, R2, and R4 pointed out, we believe that we propose a novel and well-motivated combination
 32 of multiple choice learning and model-based RL, whose major components are (i) context-conditional multi-headed
 33 dynamics model, (ii) trajectory-wise oracle loss, and (iii) adaptive planning. We also verify the effectiveness of the
 34 proposed method for improving dynamics generalization via exhaustive ablation studies on various benchmarks.

35 [R3] **Extension to other model-based RL methods.** Applying our trajectory-wise multiple choice learning scheme
 36 to model-based policy optimization methods (e.g., MBPO and Dreamer) is an interesting direction, and we think our
 37 method is naturally applicable. For example, one can consider employing our method for learning specialized policies
 38 using specialized prediction heads. We leave this as future work, but will include related discussion in the final draft.

39 [R3] **Broader Impact.** We will add more discussion related to the broader impact of our work to the final draft.

40 [R4] **Comparison to PEARL.** We emphasize that our method is evaluated using the zero-shot generalization per-
 41 formance in test environments, while PEARL is evaluated using the few-shot adaptation performance, i.e., PEARL
 42 conducts adaptation to test environments before evaluation. Our method still achieves superior sample efficiency
 43 compared to PEARL in most environments and outperforms PEARL in HalfCheetah and CrippledAnt environments
 44 even in terms of asymptotic performance. We also remark that as in PEARL, it is possible to learn a context-conditional
 45 policy using context vectors from our method (e.g., see [1]) to further improve performance.

46 [R4] **Qualitative analysis on control tasks.** This is a very interesting question. Following your suggestion, we
 47 visualize the behavior of a cheetah agent using the prediction head specialized for low-mass environments on the
 48 HalfCheetah environment with a default body mass. Interestingly, we observe that the cheetah agent moves as if it has a
 49 lightweight body, i.e., moving its limbs about very fast. We will include videos in the final draft.

50 [1] K. Lee et al. Context-aware dynamics model for generalization in model-based RL. In *ICML*, 2020.