

1 We thank all the reviewers for their thoughtful feedback. We highlight that all the reviews were positive with a few
2 specific questions, which we hope to address in our response below.

3 **Reviewer # 2**

- 4 1. Hyper-parameters : We will include an ablation study in the final version of paper.
- 5 2. How to collect “good” data for a new task : We first wish to clarify that in offline RL, the primary motivation is
6 to use existing offline datasets. As a consequence, we typically treat the dataset as fixed and given to the agent; as
7 opposed to the agent (or researcher) having the luxury of choosing the dataset. Having said that, an ideal dataset for
8 offline RL would be one that not only has high support overlap with the optimal policy, but also one that enables a
9 large hitting time, as suggested by our proposition and lower bound. Designing exploratory policies for purposes of
10 data collection is an exciting direction for future work but outside the scope of current submission.
- 11 3. Choice of NPG : We used NPG for its conceptual and implementation simplicity. A number of prior papers have
12 successfully used NPG and shown impressive results. Furthermore, our MOREL framework is modular and has
13 a clear separation between learning the P-MDP and optimizing a policy in the P-MDP. We conjecture that most
14 algorithms (e.g. PPO, SAC etc) for optimizing the policy in the P-MDP would yield similar results.

15 **Reviewer # 5**

- 16 1. Contributions of our work : To our knowledge, our work is the first to study a model-based approach to offline RL,
17 apart from Ross et al. which provided negative results for a naive algorithm. While there has been extensive work on
18 model-based RL and offline RL individually, their intersection has been explored only sparsely. As most reviewers
19 concurred, offline RL is an important learning paradigm that can expand the applicability of RL. We develop a new
20 framework for offline RL that utilizes learned models and show that it is mini-max optimal. We also demonstrate
21 state of the art experimental results.
22 Our survey of related work is extensive with 86 citations (kindly also see expanded related works in appendix). We
23 are also happy to cite and discuss any additional related work that the reviewers may point out.
- 24 2. Theoretical insights : While it is intuitively clear that if a policy does not drive too quickly towards unknown states,
25 Theorem 1 presents a *precise, quantitative* bound using *hitting times*. Corollary 3 further bounds this in terms of
26 mismatch in the support of state-action visitation distributions. In contrast, prior works only consider settings where
27 there is no support mismatch. Proposition 4 shows that the bound in Corollary 3 is best possible up to logarithmic
28 factors, demonstrating minimax optimality of MOREL.
- 29 3. Use standard errors in table : Thank you for the suggestion. We followed common practice to report standard
30 deviations, but we are happy to report standard errors if it is more appropriate in the reviewer’s opinion. Note
31 however that our claim of SOTA results in 12 out of 20 environments is based on our average scores, which remains
32 unaffected by choice of error bars. We also highlight that prior work does not report any error bars, and also tune
33 hyper-parameters on a seed-specific basis. In contrast, we use the same hyper-parameters across all seeds.
- 34 4. Proposition 4 : The value of 0.95 comes from requiring γ to satisfy certain inequalities in Lines 579 and 581 in
35 Appendix A. Since our goal is to show that the $(1 - \gamma)^{-2}$ dependence in Corollary 3 is optimal, it is okay to assume
36 that $\gamma \in [0.95, 1]$ (since, if $\gamma < 0.95$, then $(1 - \gamma)^{-2}$ is bounded by a constant).

37 **Reviewer # 6**

- 38 1. Alternate ways to penalize uncertain states : Our particular approach to penalizing unknown states enables detailed
39 theoretical analysis while also demonstrating SOTA experimental results on well studied domains that require
40 function approximation. It would make for an interesting future work to study if similar results (theoretical and/or
41 experimental) can be obtained with alternate approaches, but is outside the scope of our submission.
- 42 2. Choice of α : Note that there are two competing terms in the sub-optimality bound (Corollary 1). Decreasing α
43 decreases $(1 - \gamma)^{-2} \cdot (4\gamma R_{\max}) \cdot \alpha$, but also has the effect of decreasing the number of “known” states. This in-turn
44 reduces the hitting time, and increases $(1 - \gamma)^{-1} \cdot 2R_{\max} \cdot \mathbb{E} \left[\gamma^{T_u^{\pi^*}} \right]$ in the bound. Thus, an appropriate choice of α
45 that balances the two terms is required, and can be treated as a hyper-parameter.
- 46 3. Hyper-parameters : We will include an ablation study of hyper-parameters in the final version of the paper.
- 47 4. Comparison to more recent work : Thank you for the pointers to these very interesting papers! First, we wish
48 to highlight that these are *very recent* papers making comparisons with them difficult – especially at the time of
49 NeurIPS submission. Furthermore, these papers report results on non-standard domains compared to most prior
50 work. For example, ABM uses DeepMind control suite while BOPAH uses a different data logging policy. This
51 makes a direct comparison with published results impossible. In this submission, we used identical setups to most
52 prior papers for a fair and transparent comparison.