

## 1 **Author response for ‘When to use parametric models in reinforcement learning?’**

2 We sincerely appreciate the reviewers’ time and effort to provide useful and insightful reviews, and in this case of  
3 course especially for their helpful comments and question about our paper.

### 4 **On code and reproducibility**

5 We took care that our experiments are reproducible, and we will recheck this carefully on acceptance. We are happy to  
6 report that our main experiment has already been successfully reproduced by others, without access to our code, which  
7 is a good validation that the paper contained sufficient details. We of course appreciate that releasing code can help  
8 speed up research (and can sometimes help clarify important details), and intend to release accompanying code, and we  
9 also think it is important to make sure the paper itself contains sufficient detail to fully reproduce the results.

### 10 **On model benefits**

11 One reviewer raises interesting and insightful points about the usefulness of models that generalise better than a policy  
12 or value. We agree that this is an important aspect that is worthwhile to discuss in some detail.

13 An important, and perhaps obvious, conclusion is: it is not just important how accurate the model is (though this clearly  
14 matters), but also the model is used. We agree with the reviewer that models that generalise better than a policy or  
15 value seem especially interesting, when these are attainable. Interestingly if we then use the model in the same way as  
16 we could use replay (as opposed to, for instance, using it to plan forward for behaviour), then only accuracy does not  
17 always suffice, even in such benign settings.

18 It might be interesting to explain an experiment that we conducted, but that did not make it into the paper (yet). We set  
19 up an experiment in which the true transition dynamics were quadratic, and the model was also chosen to be quadratic.  
20 Non-surprisingly, the model very quickly learnt to match the true dynamics, essentially perfectly. To our initial surprise,  
21 the forward Dyna algorithm didn’t perform well even in this presumably benign setting—it performed far worse than  
22 the replay-based algorithm and the parameters of the value function would often diverge. We first suspected a bug but  
23 careful examination revealed this was, instead, a failure of the sort that is now discussed in Section 3. In the appendix  
24 we chose a simpler (two-state) example to concretely illustrate the more general theory around this failure, but perhaps  
25 it is useful to include this quadratic example in the paper as well, as another demonstration that even a model that  
26 generalises perfectly can fail if not used with care.

27 To be clear, of course we agree that a perfect model, if available, can help attain performance that should surpass that  
28 attained from using replay instead. We are just pointing out that using the model only to generate fictional data in the  
29 same way (and from the same states) as replay would may not be the easiest way to benefit from an imprecise model.

30 Of course, it can be hard to know a priori when a model will be accurate enough to rely on, especially for longer  
31 trajectories where compounding model errors can be a problem [cf., e.g., Talvitie, 2014, 2017, Asadi et al., 2019]. To  
32 quote Vladimir Vapnik: *one should (...) never solve a more general problem as an intermediate step* [Vapnik, 1998,  
33 Section 0.9]. Any statement of such generality comes with caveats, but it is interesting to consider in this context: when  
34 we use replay we at least know that the data is real, and we do not have to question its accuracy.

35 We agree that it is important to state our findings as clearly as possible (and appropriately scoped, to avoid hinting  
36 toward unwarranted overly general conclusions) and we intend to carefully keep polishing the writing to make the paper  
37 as clear as possible.

38 Relatedly we are also considering including further empirical results (including, though not exclusively, more at-scale  
39 results, e.g., with backward planning on Atari, as well as perhaps including the experiment with quadratic dynamics  
40 described above) to help further elucidate our main points and augment the experiments currently in the paper.

## 41 **References**

- 42 K. Asadi, D. Misra, S. Kim, and M. L. Littman. Combating the compounding-error problem with a multi-step model. *CoRR*,  
43 abs/1905.13320, 2019.
- 44 E. Talvitie. Model regularization for stable sample rollouts. In *UAI*, pages 780–789, 2014.
- 45 E. Talvitie. Self-correcting models for model-based reinforcement learning. In *Thirty-First AAAI Conference on Artificial Intelligence*,  
46 2017.
- 47 V. Vapnik. *Statistical learning theory*. Wiley, 1998.