

1 We thank the reviewers for their detailed comments and helpful suggestions. We were delighted by the enthusiasm
2 expressed by all reviewers, and the unanimous decision to place our paper above the acceptance threshold. We are
3 grateful to Reviewer 2 for the confident & positive assessment of the paper, and humbly request that Reviewers 1, 3, &
4 4 consider raising their scores from 6/10 to 7/10 if they would like this work to reach the broader NeurIPS community.
5 We will first discuss general points raised by multiple reviewers, then address reviewer-specific comments.

6 • **Concerns that REINFORCE is the only learning rule considered (R1, R2):** We apologize for this confusion—
7 we were using the term “REINFORCE” loosely, but in fact our framework can easily be applied to the family of
8 policy-gradient learning rules. In fact, two of the rules we considered, AAR and RAR, are non-REINFORCE policy-
9 gradient rules. For example, the RAR rule is derived by optimizing a different objective function compared to that for
10 REINFORCE (see SM Eq.1). In the paper, we branded them as “variants” of REINFORCE, intending to make it easier
11 for readers unfamiliar with RL language. We will clarify this distinction in our revision.

12 • **Concerns that value function based models are not considered (R1, R2):** We agree that it would be exciting to
13 explore the space of other learning models, and one of our future directions is to replace the fixed baseline that we
14 currently use, with the value function (and, thus, incorporating a TD-component) in the RF_β model. We also thank the
15 reviewers for suggesting a model comparison with variants of the Rescorla-Wagner model; this would be useful to
16 contextualize our results, and we would be happy to add this to the final paper. Meanwhile, we would like to point out
17 that while the computational cognitive science community tends to focus on TD-learning methods, a substantial number
18 of papers have proposed modeling decision-making behavior with variants of REINFORCE (e.g. Dayan & Daw (2008),
19 Kastner et al. (2019)); while Li & Daw (2011), provide support for the view that humans may use policy-gradient
20 methods instead of value prediction.

21 • **Concerns about model identifiability (R1, R3):** We agree that the identifiability of our models should be shown
22 explicitly. We will include this analysis in our revision, as well as an examination of the impact of model mismatch and
23 a broader exploration of hyperparameter space.

24

25 **Reviewer 1 :** (a) *The descriptive approach provides limited insight into how animals learn* — We apologize for
26 mischaracterizing our approach as purely descriptive; in fact, we view it as a platform for inferring the parameters
27 of normative models / testing normative hypotheses about animal learning. We would agree that our finding (that
28 negative baselines are required to account for animals’ learning trajectories) is non-intuitive, and will add supplemental
29 analyses (e.g. conditioning on incorrect choice when bias is positive) to provide more insight into how the rule affects
30 choices. (b) *Correlating weights with empirical measurements like accuracy* — This is an excellent suggestion and is
31 straightforward to show; we will certainly include supplemental figures to address this.

32 **Reviewer 2 :** (a) *Primary behavioral data to show the learning curves* — Great point, we will add learning curves
33 and other analyses to show that the inferred rules *do* indeed match behavior. (b) *Differences in learning rates for bias*
34 *and stimulus* — We should have thought of that! We will certainly discuss this in our revision, thank you for pointing
35 this out. (c) *Include additional references/lessen claims of novelty* — Thank you for the additional citations, we will
36 add these and remove the novelty claim in lines 133-4. (d) *Where is RF_β in Figure 4c/d?* — Apologies, we mistakenly
37 labeled RF_β as R+B, and RF_K as R in Fig. 4c; we will fix this in the revision.

38 **Reviewer 3 :** (a) *Evaluation of the quality of fit* — Thank you, we will add a quantification of performance in terms of
39 increase in log-likelihood on test set data. (b) *Sub-optimal behavior calls for direct validation* — Another great point,
40 we will use primary behavioral data to quantitatively validate the predictions made by the RF_β model. (c) *Statistics of*
41 *simulated vs. animal data* — Again, great point, we will include this in our revision. (d) *Outlier excluded in Fig. 4a* —
42 Apologies, this phrase was a typo and will be removed. (e) *Why is RF_β missing in Fig. 4d* — Apologies, we mistakenly
43 labeled RF_β as R+B, and RF_K as R in Fig. 4c; we will fix this in the revision.

44 **Reviewer 4 :** (a) *Clarifying the relationship to other work (ref. 19)* — Thank you for pointing this out, we will absolutely
45 make the contributions of this work clear in relation to the computational framework from [19]. (b) *Comparing full*
46 *RF_β model to RF_β without noise* — This is an interesting point, and we will include a comparison of the current RF_β fit
47 (REINFORCE with baseline, with noise) to one without a noise term in our revision. We respectfully disagree with the
48 comment that a noise term is unnecessary for the RF_K models. In fact, it was by studying the structure of the noise
49 component that we were inspired to consider the RF_β model. As was also observed, RF_β is still not perfect for the
50 animal we display in the second dataset, but through examining the structure of the retrieved noise component, we
51 may be able to suggest a better normative learning rule (which may be exactly what was suggested by the reviewer – a
52 learning rule with time-dependent baseline parameters).