

1 We would like to thank the reviewers for their time and thoughtful comments.

2 We would first like to address some clarifications:

- 3 • Using the KF does *not* make this a model-based RL method. The KF is not modeling the dynamics of the  
4 environment, rather it is estimating (with uncertainty), the local loss landscape.
- 5 • To clarify some questions on the assumptions: the assumptions of the filter depend on the model and the  
6 reward landscape, since these are the two items that influence the gradient distribution.
- 7 • Present assumptions in terms of RL.
- 8 • What do the matrices of KF represent for RL.
- 9 • TODO: Relationship to Adam: we view our method as complementary to an optimizer, like Adam. The  
10 optimizer's job is to take the gradient at a given iteration (perhaps a history of gradients as well) and perform a  
11 parameter update that improves some objective. The KF in our experiments is used to estimate the gradient at  
12 each of these iterations, not perform the update as well. Only the mean from the KF is passed to the optimizer;  
13 the variances of the KF are used only to estimate the errors.

## 14 0.1 Complexity concerns

15 Several reviewers noted potential issues in scaling to larger models. Yes, diagonal or block-diagonal approximation is  
16 an option for larger models and in these cases matrix multiplication and inversion become computationally simpler.  
17 But more generally, in control tasks there is much research suggesting that (generalized) linear models are sufficiently  
18 expressive. We chose to focus on linear models because (1) it is straightforward to ensure the assumptions hold and (2)  
19 they are sufficiently expressive for the vast majority of real-world control tasks (THIS IS TOO LARGE OF A CLAIM).

## 20 0.2 Additional experimental evaluation

21 The reviewers also noted that our experimental results section seems a bit thin and contained only a toy problem. We  
22 focused on a simple problem, as opposed to directly running the algorithm on more complicated tasks such as the  
23 MuJoCo benchmarks, because the simple LQR setup permits straight-forward evaluation and the optimal controller is  
24 known. We also focused on a generic policy gradient method (without 2nd order derivatives or a critic) because we  
25 wanted to isolate, as much as possible, the effects of the KF estimator.

26 Based on reviewer comments we ran additional experiments that compare the effects of adding a critic and also perform  
27 evaluation on more complicated benchmark tasks.

28 \* Experiments: \* Compare to additional variance reduction techniques: \* This is a valid rebuke, though we believe  
29 our method should operate in an orthogonal fashion to other approaches and did not incorporate variance reduction  
30 techniques with or without KF. \* Optimal baseline \* Linear value function baseline. \* Compare to SGD baseline. \*  
31 Additional tasks: \* Cartpole - linear. \* Walker mujoco (RBF)? \* Testing on non toy problems with a model that breaks  
32 assumptions: to test our initial work we wanted to focus on a model and task that would clearly express the inner  
33 workings of this method. A larger model, and more complicated tasks make it much more difficult to analyze anything  
34 beyond empirical performance.