

1 We would like to thank the reviewers for their detailed comments and suggestions.

2 **Scalability:** As noted by Reviewer 1, our method scales well to much larger domains. For example, our method can  
3 optimize a robust policy for a machine replacement problem with 5,000 states in only 163 seconds. We can optimize  
4 a robust policy for a 60-by-60 gridworld (3,600 states) in under two minutes. For comparison, a state-of-the-art  
5 CVaR optimization approach for MDPs with no uncertainty over the reward function takes 2 hours for a similar-sized  
6 gridworld (see [1] Section 5, last paragraph). We will add experiments demonstrating the scalability of our method  
7 to the appendix and will add an experiment where we transfer a learned reward function to a new environment (as  
8 suggested by Reviewer 1).

#### 9 **Reviewer #1**

10 > *Right now I think that the description of the contributions hides the most useful contribution.*

11 Thank you for the suggestion, we will clarify the contributions as suggested.

12 > *The robust implementation of the method (4.d) doesn't match the demonstration for one state.*

13 The reason is that Bayesian IRL does not assume demonstrator optimality, only Boltzman rationality. We used a  
14 relatively small inverse temperature resulting in reward function hypotheses that allow for occasional demonstrator  
15 errors. Using a larger inverse temperature will cause the robust policy to match all the demonstrator's actions.

16 > *The paper is fairly well written, but there is room for improvement in the paper presentation.*

17 Thank you for the detailed and constructive suggestions. We will make the notation more consistent (for example by  
18 sticking to  $\mu$  and  $w$  as much as possible) and add notation reminders throughout the paper. We will also introduce the  
19 examples as a motivation earlier in the paper.

#### 20 **Reviewer #2**

21 > *Is conditional value-at-risk the best approach for adjustable risk-sensitivity?*

22 We used CVaR because of its popularity and interpretability, but it is true that it is not always the best metric. BROIL  
23 actually works with any convex risk measure, such as EVaR or entropic risk, the only modification is that the linear  
24 program would need to be replaced by a convex optimization problem. We will make this clear in the paper.

25 > *The key issue is how to handle misspecification of the Bayesian prior.*

26 This is a good point and something we would really like to tackle in a followup work. The Bayesian statistics community  
27 has devoted a lot of effort to addressing this problem; we will add appropriate pointers.

#### 28 **Reviewer #3**

29 > *The method seems to rely heavily on the quality of prior/posterior distribution of the reward...*

30 Yes, this is true for all Bayesian methods.

31 > *The linear structure of the reward brings computational convenience, however it is hard for the reward to satisfy this  
32 structure in real-world applications.*

33 We agree that linear approximation methods (including linear regression, and linear value function approximation) have  
34 limits, but their simplicity, speed, and generally smaller data needs (bias-variance tradeoff) make them often very useful.

35 > *Any generalization of the method when we could parametrize the policy? ... deep RL or healthcare experiments?*

36 These are good suggestions. The BROIL objective is convex and nearly everywhere differentiable so it could also be  
37 used in place of expected return in a policy gradient-style approach. We judged this to be beyond the scope of this paper,  
38 but will mention this idea in the paper as an important and interesting area for future work.

#### 39 **Reviewer #4**

40 > *Assuming the samples are not exact, how do approximations in MCMC propagate onto optimization of BROIL?*

41 Thank you, this is an important question. It has been investigated in the stochastic programming community in the  
42 context of the SAA method. We will include pointers to relevant literature on this topic in the revision.

43 > *The method is presented in the context of inverse RL. Has it already been addressed within ordinary (non-inverse) RL?*

44 Several similar methods have been studied in the context of RL, but the key difference is that most RL work considers  
45 uncertain transition probabilities, while in IRL it is rewards that are uncertain. This difference has a major impact on  
46 the type of algorithms that are appropriate for the two settings. Most relevant papers in ordinary RL, which address  
47 robustness/risk aversion to *model error*, are distributionally robust MDPs (Xu & Mannor 2012), percentile optimization  
48 (Delage & Mannor 2010), and epistemic risk aversion (Eriksson & Dimitrakakis 2019).

#### 49 **References**

50 [1] Yinlam Chow, Aviv Tamar, Shie Mannor, and Marco Pavone. Risk-sensitive and robust decision-making: a cvar optimization  
51 approach. In *Advances in Neural Information Processing Systems*, pages 1522–1530, 2015.