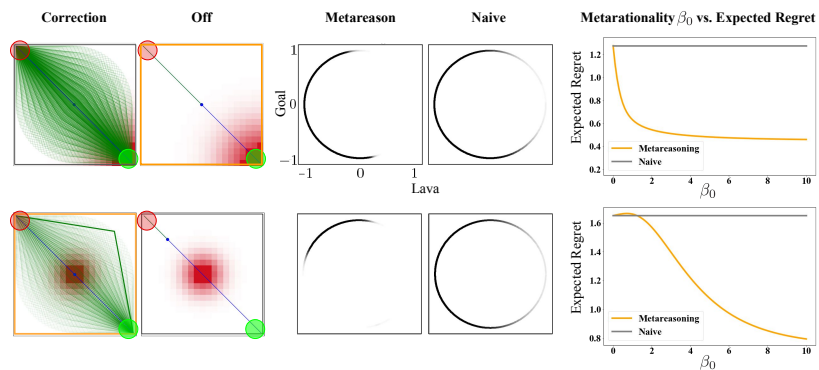


1 We thank the reviewers for their time and thoughtful feedback. They were kind to refer to the formalism as novel,  
 2 elegant, and inspiring, and to point out that the instantiations demonstrated it is fairly universal for many existing and  
 3 potentially new types of feedback. Even **R2**, our harshest critic, pointed that "I can imagine myself using this as a  
 4 resource to cite both the diversity of feedback mechanisms available and to use this formalism to develop new ones."  
 5 This is what we were hoping for! In what follows, we hope to alleviate **R2**'s main concern, and we take the opportunity  
 6 to respond to other points the reviewers brought up.

7 **Usefulness.** **R2**'s main critique is that there isn't a new method falling out of the formalism. **R4** also asks "What  
 8 are we able to do or think about that we were not able to do or think about prior to the framework?" As **R2** is  
 9 actually aware, our discussion does point to several things, including the ability to combine and actively select  
 10 the input types, but here we would like to emphasize the meta-choice. *The moment we said that there are mul-*  
 11 *tiiple types of available feedback to a person, and that we should think the person is making an implicit choice*  
 12 *within each type, it become clear that the type of feedback is itself an implicit choice* – our realization was that it  
 13 too leaks information about the reward. We actually find this to be a really compelling example of exactly what  
 14 **R2** and **R4** seem to be looking for! Now, **R2** does make a fair point that it'd be good to develop this further.

15 Unfortunately there is no way to  
 16 squeeze this in the paper (and still  
 17 explain the formalism properly). We  
 18 also really don't want it to distract  
 19 from the formalism as the main con-  
 20 tribution, which **R2** acknowledged  
 21 can already be useful in developing  
 22 new types of feedback. But we have  
 23 run some experiments with meta-  
 24 choice, and will put these in the ap-  
 25 pendix. The experiments simulate  
 26 a user choosing between corrections  
 27 and "off" when a robot was dealing  
 28 with "lava". By understanding this  
 29 as a reward rational implicit choice,  
 30 the robot is able to understand more



31 about the reward: on the bottom of the figure, if the human did a correction, and it knows the person had the "off" option  
 32 and didn't use it, that tells it about the importance of reaching the goal. The plots in the center compare the belief over  
 33 the weight on goal and lava for both naive and this "meta" inference, showing larger entropy reduction with the latter.

34 **Actual implementation.** **R1** and **R4** want to see the input types actually converted to the framework and implemented.  
 35 We want to clarify that this is what is happening in Fig.1. Those are the actual reward inferences coming from each  
 36 type, produced by our implementation (granted, in a simple domain, for illustration purposes to see how the types  
 37 compare). **R4** might be suggesting taking this "evaluation" a step further, i.e. an analysis where each feedback type is  
 38 used repeatedly. We've done experiments where the agent actively chooses which feedback is most informative which  
 39 we could add to the appendix. But we do want to (respectfully) ask the reviewer to consider that this is one of the  
 40 *many* things that would be useful to do with this framework, which is what makes the framework such a meaningful  
 41 contribution.

42 **Language.** **R4** asks why language has to result in a uniform distribution over trajectories. We apologize, it does not  
 43 have to: the formalism, as seen in eq. 1, maps choices to distributions over trajectories. Whether the distribution is  
 44 uniform or not depends on the language model. This was our mistake, we will clarify! Thank you for bringing this up.

45 **The rationality assumption.** **R1** rightfully asks whether people are actually Boltzmann-rational. While this assumption  
 46 has nice properties derived in our appendix and seems to have been useful in the works instantiating this formalism, it  
 47 is also wrong, at least when applied naively. Recent work has explored how maybe people who seem to be irrational  
 48 are actually rational, but under different assumptions. For instance, they might assume a different dynamics model, a  
 49 different observation model, or use a different planning horizon. But any such improvements in human modeling can  
 50 then translate to the formalism, now that we have all types of work under one unifying umbrella. One useful thing to  
 51 note is that an agent can potentially detect when this assumption is wrong by detecting that no reward function explains  
 52 the human's choice sufficiently well.

53 **Cost.** **R3** rightfully points out that different feedback types have different costs. When doing active learning, the agent  
 54 could trade off between information gain and user cost, or have a cost budget. We'll be sure to discuss this in the paper!  
 55 We also note that different types might be associated with different rationality parameters, which naturally affect their  
 56 informativeness.