

1 We thank the reviewers and AC for their thoughtful comments and thorough review. For this response, we have identified
2 a few common themes that have been raised by several reviewers, and address them in turn. We begin by discussing a
3 few major issues we believe are key for the reviewers’ evaluation of the main contribution of the paper:

- 4 1. Reviewers #1, #3, and #4 request a more thorough evaluation against standard RL and reward-seeking agents. A
5 comparison against a reward-seeking, exploitation-only agent is provided in Fig. 3b, showing that it performs
6 similarly to the full Active Inference (AIF) agent (but with less effective exploration, as expected). We also tested
7 agents from OpenAI’s [baselines](#) repository on the Animal-AI environment, and found that (given the same number
8 of training episodes) our agent performs considerably better than DQN and A2C, and comparable to PPO (all
9 baselines with default settings). We will include detailed comparisons in the camera-ready version of the paper.
- 10 2. Reviewer #1 urges us to describe our calculation of Eqs. 8b and 8c. For 8b, $H(s_\tau | \pi)$ is estimated sampling
11 from the transition network, and $H(s_\tau | o_\tau, \pi)$ from the encoder network (both parameterised with Gaussians, so
12 entropies can be calculated from log-variances). For the first term in 8c we sample several θ from the MC-dropouts
13 and several s_τ from the transition network; then average the entropies $H(o_\tau | s_\tau, \theta, \pi)$ (which are closed-form
14 since o_τ is Bernoulli-distributed) over the (θ, s_τ) samples. For the second term, we fix the θ and sample multiple
15 s_τ (so that, effectively, $p(o|s) = \sum_\theta p(o|s, \theta)p(\theta)$ is approximated with a single MC sample) and repeat the
16 procedure. (We also tried sampling several θ and averaging the distributions over o_τ , which is possible because o_τ
17 is Bernoulli-distributed. Although noisier, the estimator described in the paper was faster and more suitable for
18 training.) We agree with the reviewer’s statement that the entropy of the average is not the same as the average of
19 the entropies – and the difference between the two is the mutual information, which is known to be part of the EFE.
20 We will describe this calculation in detail in the appendix.
- 21 3. Reviewers #2, #3 and #4, raise concerns about the clarity of our exposition, which we group in 3 items:
 - 22 • Regarding $P(o_\tau)$: At all times it should be conditioned on π , i.e. $P(o_\tau|\pi)$. This should not have appeared in
23 lines 97 and 113 and we will update accordingly.
 - 24 • Regarding $\log P(o_\tau)$ as reward: We would like to clarify that in AIF the reward is not differentiated from
25 other types of observations. Certain (future) observations (e.g. green color in Animal-AI) are considered more
26 desirable given a task, so in practice rewards can be encoded as observations with higher prior probability
27 using $\log P(o_\tau)$. We will make this conceptual point explicit in the camera-ready version.
 - 28 • Regarding $\tilde{Q} = Q(o_\tau, s_\tau, \theta | \pi)$: The expansion is shown in Eq. 8a under the expectation, although we agree
29 it could benefit from being presented separately. We will mention it explicitly in the camera-ready version.
- 30 4. Reviewer #1 claims “scaling active inference has been done before.” We agree with the reviewer that prior work
31 has been done on this topic, but our contribution represents a technical (and qualitative) improvement over previous
32 approaches. This is achieved by: 1) estimating all summands of EFE (line 41) and, 2) for the first time successfully
33 training AIF agents on full-fledged, complex environments with visual input, multiple actions, and sparse rewards.
34 We believe this constitutes a substantial improvement over the state of the art in AIF applications.
- 35 5. Reviewers #1 and #3 claim we do not provide enough details to reproduce our results. We would like to remind
36 the reviewers we have uploaded our code to a public repository, which will be linked in the camera-ready version
37 of this paper (line 198). Additionally, for clarity we will include pseudo-code for the algorithm in the appendix.

38 In addition, we would like to address a few minor comments:

- 39 6. Reviewers #1, #2 and #3 have suggested additional references for amortised action with planning, disentanglement
40 and model-based RL. We will add these to the discussion in the camera-ready version.
- 41 7. Reviewer #2 suggests we could change the notation to be more in line with the variational inference literature.
42 Although we agree with the reviewer’s aims, given the space constraints and how much we rely on the AIF
43 literature, we believe it would make the exposition denser and the links with prior AIF literature harder to track.
44 Nonetheless, to make the paper more accessible to non-neuroscientists, in the camera-ready version we will add
45 glossary to the appendix describing in detail what each symbol and probability distribution represent.
- 46 8. Reviewer #4 argues 700 trials “seems high,” given that “one of the hallmarks of biological intelligence is few-shot
47 learning.” We agree with the reviewer, but emphasise that our agent starts ‘from scratch’ (i.e. with randomly
48 initialised networks) each run, while biological organisms are fantastically able to form good priors that generalise
49 and transfer between tasks. The extension of AIF to transfer learning remains an exciting avenue for future work.
- 50 9. The reviewers have identified a few grammatical errors, like occasional missing articles, misplaced sentences, or
51 acronyms (like ‘MC’) that should be defined more explicitly. Additionally, reviewer #2 has identified redundant
52 hyper-parameterisation (i.e. γ in Eq. 3). We will address all of these in the camera-ready version of the paper.

53 We thank again the reviewers and AC for their work, which we are sure will improve the next version of this paper. We
54 hope this response addresses the core issues raised during the review process.