

1 We thank all reviewers for their insightful comments and the time they have spent carefully reviewing the paper.
 2 Consistent among all reviewers is the comment that the paper could be improved with further experiments. In response
 3 to this, we reiterate that our aim was to provide a novel framework with a theoretically sound interpretation of RL
 4 as inference that simultaneously identifies and addresses the shortcomings of existing work while opening up new
 5 classes of algorithms within this space that others can build upon. Our experiments were designed to provide empirical
 6 evidence that our approach does not harm performance compared to state of the art, to support our theoretical claims
 7 and demonstrate acceptable performance even when the most extreme approximations are used. While we feel the
 8 submitted version already contains more than a conference paper’s worth of material, we are already running some
 9 additional experiments which, time and space permitting, we will include in the final version.

10 We will now address individual reviewer comments:

11 In response to Reviewer 1’s third comment about modelling of entire trajectories in MERLIN, the algorithms for
 12 MERLIN can all be obtained by considering the joint objective:

$$\mathcal{L}(\omega, \theta) := \mathbb{E}_{s \sim d(s)} \left[\mathbb{E}_{a \sim \pi_\theta(a|s)} \left[\frac{\hat{Q}_{\omega, \text{soft}}(h)}{\alpha} \right] \right],$$

13 where the variational distribution is $q_\theta(h) := d(s)\pi_\theta(a|s)$, the temperature constant is α and $\hat{Q}_{\omega, \text{soft}}(h)$ is the
 14 parametrised approximation for the soft action value function. The above is equivalent to max-entropy formulation in
 15 [Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review, Levine 2018] with a variational
 16 approximation for the policy. However as the variational policy trains towards the Boltzmann distribution of the
 17 soft action values with temperature α , this indirectly takes into account the dependence of the entire trajectory on
 18 the policy via $\hat{Q}_{\omega, \text{soft}}(h)$, and thus inadvertently ends up modelling entire trajectories despite grounding the MDP
 19 dynamics.

20 In response to Reviewer 2’s comment regarding the recursive definition of ϵ_ω , as we discuss in Appendix C, the
 21 definition of ϵ_ω is recursive but only if using the simple Bellman operator for the Boltzmann policy. We introduce and
 22 detail more complex operators in the set \mathbb{T} that don’t give a recursive definition in the Appendix. An example is the
 23 optimal Bellman operator, which results in a Q-learning algorithm. In Appendix F.2, we introduce another operator that
 24 recovers the optimal Bellman operator in a limit of sequences. Exploring further operators and investigating whether
 25 for any flexible Q there exist one (or many) consistent softmax temperatures $\epsilon_\omega > 0$ when using the simple Bellman
 26 operator for the Boltzmann policy is an exciting line of theoretical research for us, but one we feel is best saved for
 27 future work.

28 In response to Reviewer 2’s comment regarding comparisons to schemes where the adaptive entropy coefficient is
 29 annealed according to a schedule or optimisation scheme, as we discuss in Appendix B, our formulation has a unique
 30 Bayesian interpretation in that the entropy penalty is annealed according to the model uncertainty in the optimality of
 31 $\hat{Q}_\omega(h)$. We thank the reviewer for drawing our attention to the references [A Theory of Regularized Markov Decision
 32 Processes, Geist et al. 19] and [Soft Actor-Critic Algorithms and Applications, Haarnoja et al. 19], the former only
 33 having been published since submitting to NeurIPS, and will extend the discussion accordingly.

34 Addressing Reviewer 2’s second comment under the Quality section about function approximation for variational
 35 policies, we implied that function approximation offers the choice to obtain arbitrary rich classes of variational
 36 distributions that can, in principle, model the posterior conditional of action given the state exactly [Variational Inference
 37 with Normalizing Flows, Rezende, 2015], instead of the simpler parametrisation involving Gaussian transformations.
 38 We would like to clarify that the class of variational policies used in our experiments are the same for SAC and VIREL
 39 (both using multi dimensional independent Gaussians), thus the experiments indeed demonstrate performance gains
 40 from adaptive regularisation. We will clarify this difference in the paper.

41 Finally, we would like to thank Reviewer 3 for their careful analysis of the paper and will consider their sensible
 42 suggestion of moving details from Appendix F3 into the main paper if the NeurIPS format permits. They are correct in
 43 pointing out a reference typo in Section 5; we will update the paper accordingly.