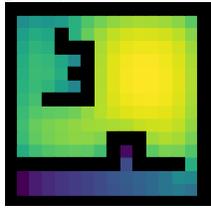


1 **R1.** Thank you for appreciating the theoretical contribution and significance. Your main concern seems the performance
 2 difference in our soft actor-critic (SAC) compared to the SAC from [20,21] on HalfCheetah. This was earlier noted by
 3 other researchers who tried to reproduce SAC results, leading to a GitHub issue which has been resolved just recently
 4 [Q. Vuong. Unable to reproduce result on HalfCheetah-v2. In *GitHub rail-berkeley/softlearning Repository*, Issue #75,
 5 2019.]. One of the SAC authors of [21] was able to reproduce low SAC performance on HalfCheetah for a specific seed
 6 setting *confirming that the effect is of statistical nature only and scientifically valid*. Details follow next. First note that
 7 the original SAC paper [20] used Mujoco v1-environments, while we used the latest v2-versions. Since then, SAC
 8 has been evaluated on v2 by the same authors as in [20]—see Figure 1 in [21]. While the performance of our SAC is
 9 comparable with the SAC from [21] on Hopper-v2, Walker2d-v2, Ant-v2 and Humanoid-v2 after 500k steps, there is a
 10 discrepancy on HalfCheetah-v2. Others obtained similar reproducing results leading to the aforementioned GitHub issue.
 11 The first figure in the first comment of the GitHub issue reports similar HalfCheetah results as we do. One SAC-paper
 12 author [21] was able to reproduce lower-performance results (see comment from May 19). The final conclusion was
 13 that this is caused by the cheetah occasionally flipping over (see comment from May 20) for specific seed settings
 14 apparently different from those used in [21]. The comment from June 24 resolves the issue: proper randomization of
 15 both environment creation and action sampling seeds in OpenAI gym leads to low performance, whereas clamping
 16 the action sampling seed to 0 yields high performance. On a minor note, the SAC results in [20,21] were obtained
 17 by averaging over 5 seeds which is not enough in Mujoco—different pairs of 5 seeds can yield significantly different
 18 results, see [43] slide 21. Therefore, we conducted experiments with 10 seeds. Hence, our evaluation is statistically
 19 more sound compared to the SAC papers [20,21]. We hope this clarification addresses your main concern.

20 **Improvements. 11)** Eq. (7) pre-states our main theoretical result in advance, but the rest of Section 4 is required to
 21 understand it—we clarified that. In short, under optimal V^* and q^* , the optimal policy π^* for the second line in Eq. (7)
 22 is given by Eq. (9). Plugging π^* back into the second line of Eq. (7) (but without the max-operator and assuming
 23 an optimal q^*) yields the third line. If q^* was replaced with a fixed q that does not depend on the next state s' , then
 24 cumulative KL regularization is recovered in which case the third line of Eq. (7) is a lower bound to the ordinary MDP
 25 formulation without logarithmic penalty. But the lower-bound statement does not hold for empowerment regularization
 26 because it *adds* intrinsic reward while KL regularization *subtracts* intrinsic reward on average. **12)** You are right. Eq. (9)
 27 can be re-formulated to yield a result similar to what you mentioned for empowerment regularization. **13)** Figure 2
 28 adopts the performance metric from [8]. We also report the plot you suggested with mean episodic reward curves in
 29 Appendix Figure 3 for all environments (we can swap them). We adjusted the paper w.r.t. **11) - 13)** accordingly.

30 **R2.** Thank you for valuing our theoretical contribution. Your main concern seems that we use a 1-step rather than
 31 a multi-step empowerment formulation. We need to stress that we optimize for *cumulative* and not *instantaneous*
 32 1-step empowerment. Cumulative 1-step empowerment yields *non-myopic* agents and *has similar properties as multi-*
 33 *step empowerment*. Note that in Figure 1 in the paper, γ was 0.6 which might evoke the impression of a myopic
 34 policy in the second plot. Below is a more illuminating example with $\gamma = 0.95$ (also with $\alpha = 0$ and $\beta = 1$).
 35  Not requiring a multi-step policy is actually a strength because a multi-step policy executes
 36 a sequence of actions (and cannot "correct" for an action when observing another state in
 37 the meanwhile). We hope this addresses your main concern and we would be happy if you
 38 shared your enthusiasm with the other reviewers (we added the new example to the paper).

39 **Improvements. Multi-Task)** We agree, empowerment could be particularly beneficial for
 40 multi-task (but this is outside the scope of the rebuttal). **Drawbacks)** We have not investigated
 41 model biases. However, empowerment specifies a particular optimization objective, and one
 42 can design reward signals that conflict with empowerment signals (e.g. negative empowerment). This could explain
 43 hindered performance—we clarified that. **Code)** We triggered the internal process for code release (in a non-academic
 44 institution, there are intellectual property regulations). **Run Time / Complexity)** Theorem 2 says for how many
 45 iterations i the value iteration needs to run to guarantee optimal values with epsilon-precision ($i \geq \log_{\gamma}(\epsilon(1-\gamma)/const)$).
 46 Proposition 3 says that one value iteration step (for a particular state s) requires an iterative Blahut-Arimoto scheme
 47 that converges at a rate of $O(1/j)$ where j is the number of "inner" iterations. Similarly to the "outer" value iteration
 48 scheme, it can be determined how many inner iterations j are required to obtain epsilon-precision, i.e. with the proof
 49 of Appendix Lemma 5: $j \geq (\beta/\epsilon)const$. The complexity of the inner Blahut-Arimoto scheme for a single state s is
 50 $O(j|S||A|)$ —see Eqs. (13,14). The overall complexity is $O(ij|S|^2|A|)$. We adjusted the paper accordingly.

51 **R3.** Thank you for your feedback. You seem concerned about the Mujoco results. While we consider the generalized
 52 MDP formulation plus theory as our main contribution, the experiments show for the first time that empowerment
 53 can improve RL in high-dimensional tasks. As you mentioned, this has been a long-standing research question not
 54 addressed before. *Empowerment leads to significant improvements in 6 tasks (most notably Ant which is amongst the*
 55 *most difficult tasks), and in the other 2 still to initial improvements* compared to the state-of-the-art SAC—see Figure 2.
 56 While we agree that more environments are always better, we already provide a suite of 8 different environments where
 57 others usually report less, e.g. 6 in SAC [20]. We adjusted the paper to address your minor details.