

# 1 Better Exploration with Optimistic Actor-Critic: Author Response

2 **All reviewers** Thanks for the feedback. As requested, we provide a plot measuring  
3 the sample efficiency gain (1) and additional ablations (Fig. 2 and 3). Also, OAC  
4 now supports deterministic policies as suggested by reviewer 1. While deterministic  
5 policies for exploration may appear surprising, deterministic OAC works because  
6 taking an action that maximises an upper bound of  $Q^\pi$  is often a better choice than  
7 taking the action that maximises the mean estimate of  $Q^\pi$ . Results are in Fig. 4.  
8 Shaded bars denote one standard deviation (runs differ due to random initial state).

9 **Reviewer 1** Thanks for the careful review. You are making several relevant points.  
10 You suggest extending OAC to support discrete, multi-modal and deterministic  
11 policies. We followed your third suggestion. We extended the scope of Proposition  
12 1 slightly to include the Wasserstein distance, deriving an OAC variant that works  
13 with deterministic policies. We report the experimental results in Figure 4, where  
14 deterministic OAC beats deterministic SAC. You also suggested (points 1 and 2 in  
15 *detailed comments*) extensions to discrete and multi-modal policies. On discrete  
16 policies, the alternative to Proposition 1 would, as you say, no longer shift the  
17 policy mean but instead constrain policy change over the probability simplex. On  
18 multimodal-policies, extending OAC to Gaussian mixtures akin to Actor-Expert  
19 (Lim, 2018) would imply taking the KL divergence between mixtures. We agree  
20 that these extensions are interesting, but they would be hard to pack into a single  
21 submission. We will discuss them in the future work section and also relate our  
22 algorithm to Actor-Expert (Lim, 2018). In point 3, you suggested making an  
23 optimistic variant of TD3 or DDPG. We agree this would be informative, but we  
24 had a limited computation budget and chose SAC because of its performance on  
25 Humanoid. On your point 4, as you say, our results may not be groundbreaking but  
26 the difference is statistically significant and a step forward (Fig. 1). Also, thanks  
27 for flagging the TD3 results. There was a problem in our setup and the results are  
28 now 5K on Ant after 2.5M steps. Concerning the *small issues* part of the review,  
29 we will clarify description of TD3, fix the misnumbered equation and discuss the  
30 cost of computing the additional gradient (it is very small in practice).

31 **Reviewer 2** We appreciate the kind words. You are right when you say that  
32 OAC still needs many environmental interactions. However, using OAC vs SAC  
33 does make a meaningful difference. On Humanoid, OAC obtains a policy of same  
34 quality in 0.52M steps vs 1M for SAC (Figure 1). We agree that improving sample  
35 efficiency remains a challenge and hope that OAC paves the way for even better  
36 methods.

37 **Reviewer 3** Thanks for the feedback. In *detailed comments*, you ask why a spurious  
38 maximum of the lower bound leads to a policy with small covariance. Intuitively,  
39 this is because the actor finds a probability distribution that greedily maximises the  
40 critic lower bound. But a distribution that maximises a function is a point mass at  
41 the maximum of that function. Formally, as actor iteration progresses, the covariance  $\Sigma$  can be modelled as  $e^{Ht}$ , where  
42  $H$  is the second order term in the Taylor expansion of the critic around the policy mean and  $t$  is the iteration count. Near  
43 a maximum,  $H$  is negative definite and we have  $\Sigma \propto e^{Ht} \rightarrow 0$  as  $t \rightarrow \infty$ . We will include an extension of this argument  
44 in the paper. Your second point concerns how our Gaussian exploration policy avoids directional uninformedness. This  
45 is best seen in the figure on page 5 of the paper. While the exploration policy  $\pi_E$  is symmetric around its own mean, it is  
46 not symmetric around the mean of the target policy  $\pi_T$ . We will make this clearer. Also, you requested a measurement  
47 of the directionality. We provide it in Fig. 3, which tracks the absolute magnitude in the difference between the mean  
48 of the exploration policy and target policy. We also performed an ablation for optimism, shown in Fig. 2. The figure  
49 shows a sweet spot (the optimism value  $\beta_{UB} = 4.36$  we used in the submission). About proposition 1, we will motivate  
50  $\Sigma_E = \Sigma_T$  more clearly. We will also expand the justification for this near line 450 of Appendix A. Also, you propose  
51 using  $\max(\hat{Q}_{LB}^1, \hat{Q}_{LB}^2)$  as the UCB. We in fact already do, i.e.  $\max(\hat{Q}_{LB}^1, \hat{Q}_{LB}^2) = \mu_Q + \sigma_Q$ , using notations from lines  
52 154-158 and Appendix B. We will make this clearer. Also, as requested, we provide a plot measuring the number of  
53 steps to reach a given performance (Fig. 1). On your minor comments, we meant actor-critic in line 73. We will fix this  
54 as well as the typo, the misnumbered equation and the format of references.

## 55 References

56 S. Lim, J. Ajin L. Le, Y. Pan, M. White *Actor-Expert: A Framework for using Action-Value Methods in Continuous Action Spaces*

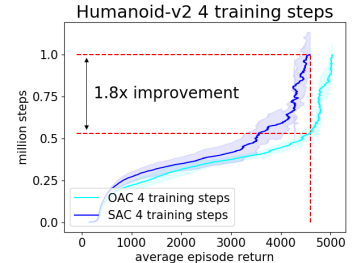


Fig. 1: Sample efficiency

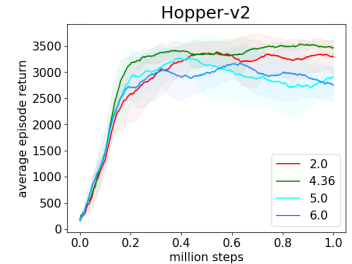


Fig. 2: Optimism ablation ( $\beta_{UB}$ )

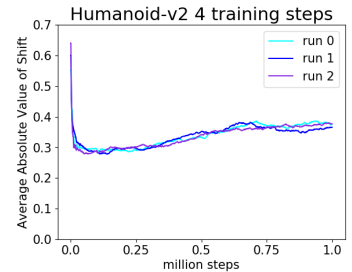


Fig. 3: Magnitude of shift.

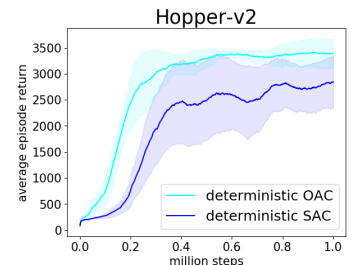


Fig. 4: Deterministic OAC