

1 We thank the reviewers for their interest in our work and their helpful comments. Please find our response below.

2 **Comments relevant to all reviewers:**

3 - In the three time scale procedure, faster time scales view slower time scales as static. This is why the fastest time scale  
4 is essentially solving a supervised learning problem over two static networks.

5 - The slowest time scale (delayed actor) is required both theoretically as well as empirically. Without it, convergence is  
6 not guaranteed and the algorithm becomes unstable.

7 **Reviewer 1:**

8 - You are correct in pointing out that an on-policy version of Algorithm 1 is not ensured to converge. DPO is an  
9 off-policy actor-critic framework which requires that all state action pairs are visited "enough" in order to ensure  
10 convergence, which is a theoretical assumption in various off-policy algorithms. We achieve this, similar to DQN,  
11 DDPG and TD3, by keeping an exploration strategy which does not decay to zero. We will emphasize this to avoid  
12 confusion in the paper.

13 - Your observation is correct, DPO and GAC are not perfectly aligned. DPO requires optimization in the distribution  
14 space, while GAC is a practical approximation of DPO in which optimization occurs in the space of parameters of the  
15 generative model. The DPO framework is a fundamental framework which can be extended in a similar way as Policy  
16 Gradient methods to bridge the gap between DPO and GAC.

17 - GANs and VAEs are definitely a valid choice for representing the policy, yet they have some pitfalls [1]. GANs pose  
18 the problem of learning a generative model and solving a two-player zero-sum game. This form of learning in itself  
19 is often unstable (resulting in mode-collapse) and still lacks theoretical guarantees and stability assurances. VAEs  
20 minimize the KL distance, as opposed to the p-Wasserstein distance, which has its own benefits. The quantile approach  
21 overcomes these issues by directly minimizing the p-Wasserstein distance using the quantile regression loss.

22 **Reviewer 2:**

23 - This is a good point. In practice, the three time-scale requirement is implemented using different learning rates so that  
24 the various elements converge at different rates. We follow a similar implementation method as in other actor-critic  
25 approaches, which are based on two timescales.

26 - *Derivation of policy distribution update:*

27 What we ultimately wish to have is an update similar to that of Policy Iteration (or more specifically, Approximate  
28 Policy Iteration), which conservatively updates the policy given a target policy  $\pi'$  as:

$$\pi_{k+1}(a|s) = (1 - \alpha_k)\pi_k(a|s) + \alpha_k\pi'(a|s).$$

29 Policy Iteration schemes use the target policy  $\pi'(a|s) \in \arg \max_{a \in \mathcal{A}} r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)v^{\pi_k}(s')$  in the exact  
30 case, or the  $\epsilon$ -greedy target policy in the approximate case. Nevertheless, finding the  $\arg \max$  is itself a hard problem in  
31 non-convex continuous regimes. Since finding the greedy action is complicated, it is reasonable to instead define the  
32 target policy (i.e.,  $\pi'$ ) as a distribution over all improving actions. We denote this target policy by  $\pi'(a|s) = \mathcal{D}_{I^{\pi_k}}^{\pi}(a|s)$ .  
33 Finally, we take a gradient based approach using a distance metric over policies,  $d$ , yielding the DPO update rule

$$\pi_{k+1} = \Gamma(\pi_k - \alpha_k \nabla_{\pi} d(\mathcal{D}_{I^{\pi_k}}^{\pi}, \pi) |_{\pi=\pi_k}).$$

34 **Reviewer 3:** Thank you for pointing out some confusing explanations, we will make sure to clarify them in the paper.

35 - In Fig. 1a the intention is to compare optimization w.r.t. the policy itself (i.e.,  $\nabla_{\pi} v^{\pi}$ ) to optimization w.r.t. the  
36 parametrization  $\nabla_{\theta} v^{\pi_{\theta}}$  (e.g., for Delta distributions  $\theta$  represents the action, and for Gaussian distributions  $\theta$  represents  
37 the mean  $\mu$ ). The former is what classical algorithms such as CPI (Kakade and Langford 2002) require, whereas the  
38 latter is what occurs in the standard policy gradient approaches. In Fig. 1a, the left ( $\Pi$  space) represents the ideal  
39 approach, whereas the right ( $\Theta$  space) represents the sub-optimality which occurs when encountering a non-convex  
40 action-value function and optimizing with respect to the parametric distribution parameters (e.g., action).

41 - Regarding convexity, our intent was to show that the set  $\Theta$  is not convex in a probabilistic sense. The set  $\Pi$  is the set  
42 of all probability distributions, whereas  $\Theta$  is the span of probabilities distributions that  $\pi_{\theta}$  can represent. Gaussian or  
43 Delta distributions are limited to their set, and thus can't ensure convergence to a global extrema. More specifically,  
44  $\alpha\delta_{\mu_1} + (1 - \alpha)\delta_{\mu_2}$  means to play action  $a_1 = \mu_1$  with probability  $\alpha$  and  $a_2 = \mu_2$  otherwise, but this distribution is not  
45 a Delta distribution, and therefore is not contained in  $\Theta$  - the set of all Delta distribution functions. We will make this  
46 notation clear in the paper.

47 - Thank you for noting the mistake in lines 99 and 442 - we've updated the paper.

48 [1] Georg Ostrovski, Will Dabney, and Remi Munos. Autoregressive quantile networks for generativemodeling. *arXiv*  
49 *preprint arXiv:1806.05575*, 2018.