

1 We would like to thank the reviewers for their thorough evaluations and for bringing to our attention some missing
 2 citations and typos, these will be corrected in the updated manuscript. We answer here specific questions raised by the
 3 reviewers and present requested additional experiments.

4 **Additional MuJoCo experiments (R1, R2, R4).** We focus on GAIL because AIRL claims to perform on par with GAIL
 5 on MuJoCo. We present in Figure 2 below additional results on the MuJoCo envs (as well as the additional Ant env)
 6 where GAIL has been re-tuned for further improvement and SQIL and ASAF-1 have been added. We see that even with
 7 careful tuning GAIL is outperformed by our method and that SQIL’s instability is exacerbated on MuJoCo.

8 We also ran ASAF-1 on the Ant-v2 MuJoCo environment using various sets
 9 of 25 demonstrations (*as requested by R2*). These demonstrations were
 10 generated from a Soft Actor-Critic agent at various levels of performance
 11 during its training. Since at low-levels of performance the variance of
 12 episode return is high, we filtered collected demonstrations to lie in the
 13 targeted range of performance (e.g. return in [800, 1200] for the 1K set).
 14 Results in Figure 1 show that our algorithm succeeds at learning a policy
 15 that closely emulates various demonstrators (even when non-optimal).

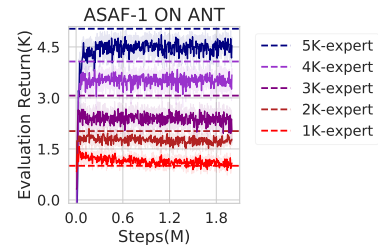


Figure 1: ASAF-1 on Ant-v2. Colors are 1K, 2K, 3K, 4K, 5K expert’s performance.

16 **Comparison of our methods vs Finn et al. [1] and Fu et al. [2] and**
 17 **Behavioral Cloning (BC) (R1, R3).** Our work builds on [1], yet its novelty
 18 is to explicitly express the probability of a trajectory in terms of the policy in
 19 order to directly learn this latter when training the discriminator. In contrast, [2] considers a transition-wise discriminator
 20 with un-normalized probabilities which makes it closer to ASQF (Appendix B) than to ASAF-1. Additionally, AIRL
 21 [2] minimizes the Kullback-Leiber Divergence [3] between occupancy measures whereas ASAF minimizes the Jensen-
 22 Shannon Divergence between trajectories likelihood. Finally, BC uses the loss function from supervised learning
 23 (classification or regression) to match expert’s actions given expert’s states and suffers from compounding error due
 24 to co-variate shift [4] since it only learns on the expert state-action visitations (demonstrations) without environment
 25 interaction. Contrarily, ASAF-1 uses the binary cross entropy loss in Eq. (13) and does not suffer from compounding
 26 error as it learns on both generated and expert’s trajectories.

27 **Can windowed approach be used for GAIL and AIRL? (R1).** GAIL’s and AIRL’s discriminators are updated based on
 28 transitions so that a step-wise reward can be learned and used in RL loops. Therefore, the windowed approach, which
 29 suits trajectory-wise formulations, could not be applied without major modification of their formulations.

30 **Reward Acquisition from ASAF (R3).** Although ASAF does not explicitly acquire reward during training, we can
 31 retrieve step-wise soft advantages $\log \pi(a|s)$ from learned agent’s policy $\pi(a|s)$ which can be used as a reward
 32 function [2, 5].

33 **Training Time (R4).** Due to lack of room we cannot add here the equivalent of Figure 2 with wall clock time as x-axis
 34 but we will add it to the updated manuscript. Our observation is that ASAF-1 is always fastest to learn, e.g., 361.2s
 35 (ASAF-1), 473.1s (SQIL), 1561.0s (ASAF), 1763.6s (GAIL), 2079.1s (ASAF-w) were taken to reach a performance of
 36 2000 in Hopper environment. Note however that reports of performance w.r.t wall-clock time should always be taken
 37 with a grain of salt as they are greatly influenced by hyper-parameters and implementation details.

38 **No entropy regularization in loss (R3).** We learn in the softmax policy class of Eq. (2) since it contains the expert’s
 39 policy (See Section 3.2) but optimize Eq. (13) without entropy regularization as was done in the GAIL paper.

40 **Additional Concerns. (To R3)** We will update our related works with your recommendations. **(To R4)** It seems that our
 41 supplementary files have been correctly uploaded since R1 was able to read through our code. We are sorry that you
 42 couldn’t access it. Unfortunately, we are not allowed to put an external link in this rebuttal. We suggest that you contact
 43 the AC about this issue.

44 **References.** [1] Finn et al., "A connection between generative ...," (2016) [2] Fu et al., "Learning robust rewards ...,"
 45 (2017) [3] Ghasemipour et al., "A divergence minimization perspective ...," (2019) [4] Ross and Bagnell, "Efficient
 reductions for imitation learning," (2010) [5] Schulman et al., "High-dimensional continuous control ...," (2015)

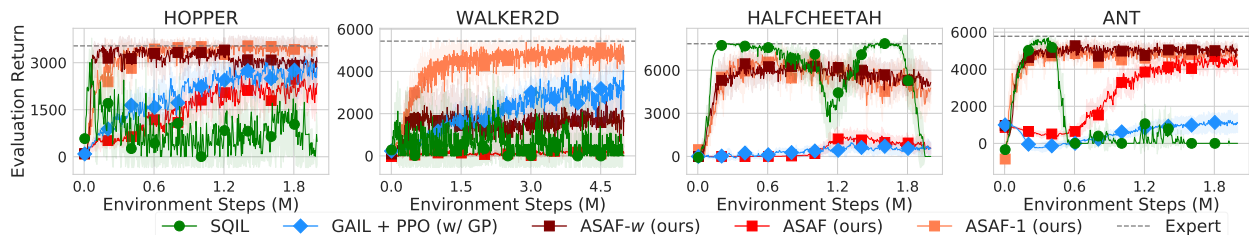


Figure 2: Results on Mujoco environments with improved tuning for GAIL, added SQIL and ASAF-1 and env (Ant)