

1 Thank reviewers for the comments. Please find our responses below, *with reference indices consistent with the paper.*

2 **To reviewer 3. Q3-1: Experimental choices.** Imitation learning (IL)  
 3 aims to match the state-action distributions between the learner and the  
 4 expert, rather than the goal feature(s). The expert state-action distributions  
 5 from the data are high-dimensional and stochastic. Taking CartPole as  
 6 an example, Fig. 1 (right) shows the expert (stochastic) policy on two  
 7 sample states, where the states are from a 4-dimensional continuous space  
 8 (see Appendix B in [18]). Moreover, Fig. 1 (left) shows that in the expert  
 9 demonstrations, the angle (feature) is NOT simply a constant (over 250  
 10 trajectories in CartPole). In addition, for fair comparisons, our evaluation

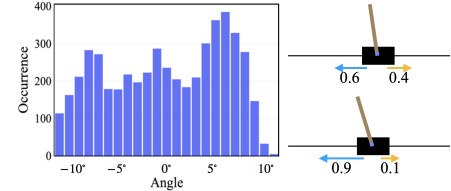


Figure 1: Angle dynamics (left) and policy (right) in CartPole from expert data.

11 tasks/settings are consistent with SOTA [13,18,20]. **Q3-2: Sample efficiency.** Results on sample efficiency are  
 12 presented in Fig. 5 in the paper, where  $f$ -GAIL outperforms baselines over different sample sizes (Sec 4.2). Moreover,  
 13 it is true that  $f$ -GAIL can more accurately estimate  $f$ -divergence from a limited number of samples. Fig. 2 below shows  
 14 that given a task, the learned  $f^*$  functions are consistent with different sample sizes. In fact, the choice of  $f$ -divergence  
 15 matters. The better divergence estimation accuracy enables  $f$ -GAIL to examine and compare  $f$ -divergence choices,  
 16 which is why  $f$ -GAIL consistently outperforms baselines.

17 **Q3-3: Meaning of the “best”  $f$ -divergence.** Our  $f$ -  
 18 GAIL is defined as a minimax optimization problem in  
 19 eq.(5) in the paper. The best  $f$ -divergence is searched in  
 20 the “max” inner-loop given the current learned policy  
 21  $\pi$  learned from the “min” outer-loop, eventually leading  
 22 to a stable solution of  $(\pi, f^*)$ . **Q3-4: The optimality  
 23 depends on the divergence and context?** The notion of  
 24 optimality exists and depends on the expert demonstration  
 25 data, rather than the divergence, namely, the optimality

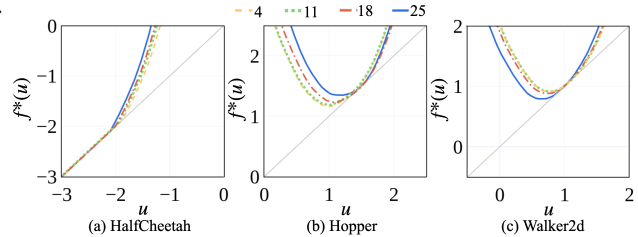


Figure 2: The learned  $f^*_\phi(u)$  with different sample sizes

26 refers to the smallest discrepancy of behavior distributions (in state-action pairs) between the learner and the expert.  
 27 Given an expert demonstration dataset, a better divergence can measure the discrepancy more precisely than other  
 28 divergences, thus enable training a learner with closer behaviors to the expert. In the example of whether mode-seeking  
 29 or -covering makes sense, it depends on the expert demonstration data (NOT the context), i.e., whether the expert was  
 30 performing mode-seeking or -covering when generating the data. **Q3-5: Meaning of the learned divergence? The  
 31 variance of the learned divergence?** It is nontrivial to find an analytical close-form function to express the learned  
 32  $f^*$ , due to the huge convex function space with  $f(1) = 0$ . We leave this as our future work. Fig. 2 above shows that  
 33 given a task, the learned  $f^*$  functions are consistent (small variance) for different sample sizes. **Q3-6: Comparison  
 34 with BC.** We agree that BC minimizes the policy KL divergence as what we noted in Sec. 4 (line 200). We included  
 35 BC as a baseline for completeness, namely, a comprehensive comparison with SOTA. **Q3-7: Notation of  $\mathcal{P}$  in line 62.**  
 36 Our notation represents a probability distribution of transitioning from  $(s, a)$  to a next state  $s'$ , thus the outcome is in  
 37  $[0, 1]$ . It is consistent with the literature, e.g., Sec. 2 in [Yu et al. arXiv:1909.09314].

38 **To reviewer 2. Q2-1: Necessity and sufficiency of two constraints.** The  $f$ -  
 39 divergence definition requires the generator  $f$  function to be convex and  $f(1) = 0$   
 40 [11,23,24]. Convex and zero-gap constraints are necessary and sufficient conditions to  
 41 guarantee an  $f$ -divergence, based on  $f^{**} = f$  (see §3.3.2 in [9]) for convex functions,  
 42 i.e.,  $f(1) = f^{**}(1) = \max_u \{u - f^*(u)\} = 0$ . **Q2-2: Implementation details.** We

43 used 5 random seeds with mean and variance calculated over 50 trajectories (see Sec. 4 and Appendix B). These  
 44 settings are consistent with SOTA [11,23,24]. **Q2-3: Baselines.** All baselines were implemented with their original  
 45 models [13,18,20], rather than  $f^*(T(s, a))$ . In fact, as shown in Tab. 1 (with GAIL as the original model and  $GAIL_f$   
 46 as  $f^*(T(s, a))$ ), the baseline results are similar, when implemented using original models vs  $f^*(T(s, a))$ . **Q2-4:** The  
 47 input state distributions were sampled from expert demonstrations. **Q2-5: Divergence evaluation.** Following your  
 48 suggestion, Fig. 3 below shows the training curve of  $f$ -divergence wrt. epochs where it converges to less than 0.02 for  
 49 HalfCheetah after 450 epochs. Similar results were observed in other tasks.

50 **To reviewer 1. Q1-1: Novelty:** We are the *first* to model imitation learning with a learnable  
 51  $f$ -divergence measure (using the proposed  $f^*$ -network), rather than a predefined divergence,  
 52 which yields better learner policies than the literature on GAIL [13,18,20]. **Q1-2: More  
 53 complex tasks:** We evaluated  $f$ -GAIL on tasks consistent with SOTA [13,18,20], including  
 54 Humanoid, with the high state dimension of 376. We plan to evaluate  $f$ -GAIL on more  
 55 complex tasks, e.g., Simitate [Memmesheimer et al. arXiv:1905.06002].

56 **To reviewer 5. Training objective:** In our Alg 1, all three networks are trained with the  
 57 same objective in eq.(5), using *adversarial training*. The updating gradients ( $\nabla_\omega$  and  $\nabla_\phi$ )  
 58 are obtained by taking the derivative of eq.(5) wrt.  $\omega$  and  $\phi$ , respectively, while fixing  $\pi_\theta$ .  
 59 The objective for policy  $\pi_\theta$  is the same as eq.(5) with  $T_\omega$  and  $f^*_\phi$  fixed. We will add such details in the final paper.

Table 1: Baseline performances in HalfCheetah.

Datasize	GAIL	$GAIL_f$
4	$4047 \pm 344$	$4055 \pm 257$
25	$4340 \pm 185$	$4472 \pm 166$

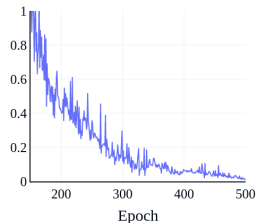


Figure 3:  $f$ -divergence curve in HalfCheetah.