

---

# A Ranking Game for Imitation Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

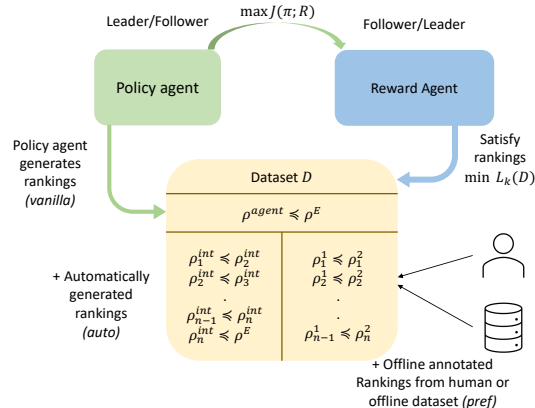
1 We propose a new framework for imitation learning—treating imitation as a *two-*  
2 *player ranking-based game* between a policy and a reward. In this game, the reward  
3 agent learns to satisfy pairwise performance rankings between behaviors, while the  
4 policy agent learns to maximize this reward. In imitation learning, near-optimal  
5 expert data can be difficult to obtain, and even in the limit of infinite data cannot  
6 imply a total ordering over trajectories as preferences can. On the other hand,  
7 learning from preferences alone is challenging as a large number of preferences  
8 are required to infer a high-dimensional reward function, though preference data is  
9 typically much easier to collect than expert demonstrations. The classical inverse  
10 reinforcement learning (IRL) formulation learns from expert demonstrations but  
11 provides no mechanism to incorporate learning from offline preferences and vice  
12 versa. We instantiate the proposed ranking-game framework with a novel ranking  
13 loss giving an algorithm that can simultaneously learn from expert demonstrations  
14 and preferences, gaining the advantages of both modalities. Our experiments show  
15 that the proposed method achieves state-of-the-art sample efficiency and can solve  
16 previously unsolvable tasks in the Learning from Observation (LfO) setting.

## 17 1 Introduction

18 Reinforcement learning relies on environmental reward feedback to learn meaningful behaviors.  
19 Reward specification is a hard problem [39], thus motivating imitation learning (IL) as a technique  
20 to bypass reward specification and learn from expert data, often via Inverse Reinforcement Learning  
21 (IRL) techniques. Learning from expert **observations** (imitation learning) alone can require efficient  
22 exploration when the expert actions are unavailable as in LfO [36]. Incorporating preferences over  
23 potentially suboptimal trajectories for reward learning can help reduce the exploration burden by  
24 regularizing the reward function and providing effective guidance for policy optimization. Previous  
25 literature in learning from preferences either assumes no environment interaction [10, 9] or assumes  
26 an active query framework with a restricted reward class [47]. The classical IRL formulation suffers  
27 from two issues: (1) Learning from expert demonstrations and learning from preferences/rankings  
28 provide complementary advantages for increasing learning efficiency [30, 47]; however, existing  
29 IRL methods that learn from expert demonstrations provide no mechanisms to incorporate offline  
30 preferences and vice versa. (2) Optimization is difficult, making learning sample inefficient [5, 28]  
31 due to the adversarial min-max game.

32 Our primary contribution is an algorithmic framework casting imitation learning as a rank-  
33 ing game which addresses both of the above issues in IRL. This framework treats imi-  
34 tation as a ranking game between two agents: a reward agent and a policy agent—the  
35 reward agent learns to satisfy pairwise performance rankings between different *behaviors*

36 represented as state-action or state visitations, while the policy agent maximizes its per-  
 37 formance under the learned reward function. The ranking game is detailed in Figure 1  
 38 and is specified by three components: (1) The dataset of pairwise behavior rankings, (2) A  
 39 dataset of pairwise behavior rankings, (2) A ranking loss function, and (3) An optimization  
 40 strategy. This game encompasses a large subset of both inverse reinforcement learning (IRL)  
 41 methods and methods which learn from suboptimal offline preferences. Popular IRL methods  
 42 such as GAIL, AIRL,  $f$ -MAX [28, 22, 34] are instantiations of this ranking game in which  
 43 rankings are given only between the learning agent and the expert, and a gradient descent  
 44 (GDA) optimization strategy is used with a ranking loss that maximizes the performance  
 45 gap between the behavior rankings.  
 46  
 47  
 48  
 49  
 50  
 51



52 The ranking loss used by the prior IRL  
 53 approaches is specific to the comparison of  
 54 optimal (expert) vs. suboptimal (agent) data,  
 55 and precludes incorporation of comparisons  
 56 among suboptimal behaviors. In this work, we  
 57 instantiate the ranking game by proposing a new  
 58 ranking loss ( $L_k$ ) that facilitates incorporation  
 59 of [rankings over suboptimal trajectories](#) for  
 60 reward learning. Our theoretical analysis reveals that the proposed ranking loss results in a bounded  
 61 performance gap with the expert that depends on a controllable hyperparameter. Our ranking loss  
 62 can also ease policy optimization by supporting data augmentation to make the reward landscape  
 63 smooth and allowing control over the learned reward scale. Finally, viewing our ranking game in the  
 64 Stackelberg game framework (see Section 3)—an efficient setup for solving general-sum games—we  
 65 obtain two algorithms with complementary benefits in non-stationary environments depending on  
 66 which agent is set to be the leader.

Figure 1: rank-game: The Policy agent maximizes the reward function by interacting with the environment. The Reward agent satisfies a set of behavior rankings obtained from various sources: generated by the policy agent (vanilla), automatically generated (auto), or offline annotated rankings obtained from a human or offline dataset (pref). Treating this game in the Stackelberg framework leads to either Policy being a leader and Reward being a follower, or vice-versa.

67 In summary, this paper formulates a new framework rank-game for imitation learning that allows  
 68 us to view learning from preferences and demonstrations under a unified perspective. We instantiate  
 69 the framework with a principled ranking loss that can naturally incorporate rankings provided by di-  
 70 verse sources. Finally, by incorporating additional rankings—auto-generated or offline—our method:  
 71 (a) outperforms state-of-the-art methods for imitation learning in several MuJoCo simulated domains  
 72 by a significant margin and (b) solves complex tasks like imitating to reorient a pen with dextrous ma-  
 73 nipulation using only a few observation trajectories that none of the previous LfO baselines can solve.

## 74 2 Related Work

75 Imitation learning methods are broadly divided into two categories: Behavioral cloning [48, 54] and  
 76 Inverse Reinforcement Learning (IRL) [44, 1, 72, 18, 20, 28, 22]. Our work focuses on developing a  
 77 new framework in the setting of IRL through the lens of ranking. Table 1 shows a comparison of the  
 78 proposed rank-game method to prior works.

79 **Classical Imitation Game for IRL:** The classical imitation game for IRL aims to solve the  
 80 adversarial *min-max* problem of finding a policy that minimizes the worst-case performance gap  
 81 between the agent and the expert. A number of previous works [22, 60, 34] have focused on  
 82 analyzing the properties of this *min-max* game and its relation to divergence minimization. Under  
 83 some additional regularization, this *min-max* objective can be understood as minimizing a certain  
 84  $f$ -divergence [28, 22, 34] between the agent and expert state-action visitation. More recently, [60]  
 85 showed that all forms of imitation learning (BC and IRL) can be understood as performing moment  
 86 matching under differing assumptions. In this work, we present a new perspective on imitation

IL Method	Offline Preferences	Expert Data	Ranking Loss	Reward Function	Active Human Query
MaxEntIRL, AdRIL, GAN-GCL, GAIL, $f$ -MAX, AIRL	✗	LfD	supremum	non-linear	✗
BCO, GAfO, DACfO, OPOLO, $f$ -IRL	✗	LfO	supremum	non-linear	✗
TREX, DREX	✓	✗	Bradley-Terry	non-linear	✗
BREX	✓	✗	Bradley-Terry	linear	✗
DemPref	✓	LfO/LfD	Bradley-Terry	linear	✓
Ibarz et al[30]	✓	LfD	Bradley-Terry	non-linear	✓
rank-game	✓	LfO/LfD	$L_k$	non-linear	✗

Table 1: A summary of IL methods demonstrating the data modalities they can handle (expert data and/or preferences), the ranking-loss functions they use, the assumptions they make on reward function, and whether they require availability of an external agent to provide preferences during training. We highlight whether a method enables LfD, LfO, or both when it is able to incorporate expert data.

87 in which the reward function is learned using a dataset of behavior comparisons, generalizing  
88 previous IRL methods that learn from expert demonstrations and additionally giving the flexibility  
89 to incorporate rankings over suboptimal behaviors.

90 **Learning from Preferences and Suboptimal Data:** Learning from preferences and suboptimal data  
91 is important when expert data is limited or hard to obtain. Preferences [3, 65, 55, 14, 47] have the  
92 advantage of providing guidance in situations expert might not get into, and in the limit provides  
93 full ordering over trajectories which expert data cannot. A previous line of work [10, 11, 9, 13] has  
94 studied this setting and demonstrated that offline rankings over suboptimal behaviors can be effectively  
95 leveraged to learn a reward function. [14, 47, 30] studied the question of learning from preferences  
96 in the setting when a human is available to provide online preferences<sup>1</sup> (active queries), while [47]  
97 additionally assumed the reward to be linear in known features. Our work makes no such assumptions  
98 and allows for integrating offline preferences and expert demonstrations under a common framework.

99 **Learning from Observation (LfO):** LfO is the problem setting of learning from expert observations.  
100 This is typically more challenging than the traditional learning from demonstration setting (LfD),  
101 because actions taken by the expert are unavailable. LfO is broadly formulated using two objectives:  
102 state-next state marginal matching [63, 71, 58] and direct state marginal matching [45, 43]. Some  
103 prior works [61, 67, 16] approach LfO by inferring expert actions through a learned inverse dynamics  
104 model. These methods assume injective dynamics and suffer from compounding errors when the  
105 policy is deployed. A recently proposed method OPOLO [71] derives an upper bound for the LfO  
106 objective which enables it to utilize off-policy data and increase sample efficiency. Our method  
107 outperforms baselines including OPOLO, by a significant margin.

### 108 3 Background

109 We consider a learning agent in a Markov Decision Process (MDP) [49, 59] which can be defined  
110 as a tuple:  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \rho_0)$ , where  $\mathcal{S}$  and  $\mathcal{A}$  are the state and action spaces;  $P$  is the state  
111 transition probability function, with  $P(s'|s, a)$  indicating the probability of transitioning from  $s$  to  $s'$   
112 when taking action  $a$ ;  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function bounded in  $[0, R_{max}]$ ; We consider  
113 MDPs with infinite horizon, with the discount factor  $\gamma \in [0, 1]$ , though our results extend to finite  
114 horizons as well;  $\rho_0$  is the initial state distribution. We use  $\Pi$  and  $\mathcal{R}$  to denote the space of policies  
115 and reward functions respectively. A reinforcement learning agent aims to find a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$   
116 that maximizes its expected return,  $J(R; \pi) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \rho^\pi(s,a)} [R(s, a)]$ , where  $\rho^\pi(s, a)$  is the  
117 stationary state-action distribution induced by  $\pi$ . In imitation learning, we are provided with samples  
118 from the state-action visitation of the expert  $\rho^{\pi^E}(s, a)$  but the reward function of the expert is  
119 unknown. We will use  $\rho^E(s, a)$  as a shorthand for  $\rho^{\pi^E}(s, a)$ .

120 **Classical Imitation Learning:** The goal of imitation learning is to close the imitation gap  $J(R; \pi^E) -$   
121  $J(R; \pi)$  defined with respect to the unknown expert reward function  $R$ . Several prior works [28, 60,  
122 38, 45] tackle this problem by minimizing the imitation gap on all possible reward hypotheses. This

<sup>1</sup>We will use preferences and ranking interchangeably

123 leads to a zero-sum (min-max) game formulation of imitation learning in which a policy is optimized  
 124 with respect to the reward function that induces the largest imitation gap:

$$\text{imit-game}(\pi) = \arg \min_{\pi \in \Pi} \sup_{f \in \mathcal{R}} \mathbb{E}_{\rho^E(s,a)}[f(s,a)] - \mathbb{E}_{\rho^\pi(s,a)}[f(s,a)]. \quad (1)$$

125 Here, the imitation gap is upper bounded as follows ( $\forall \pi$ ):

$$J(R; \pi^E) - J(R; \pi) \leq \sup_{f \in \mathcal{R}} \mathbb{E}_{\rho^E(s,a)}[f(s,a)] - \mathbb{E}_{\rho^\pi(s,a)}[f(s,a)]. \quad (2)$$

126 Note that, when the performance gap is maximized between the expert  $\pi^E$  and the agent  $\pi$ , we can  
 127 observe that the worst-case reward function  $f_\pi$  induces a ranking between policy behaviors based  
 128 on their performance:  $\rho^E \succeq \rho^\pi := \mathbb{E}_{\rho^E(s,a)}[f_\pi(s,a)] \geq \mathbb{E}_{\rho^\pi(s,a)}[f_\pi(s,a)]$ ,  $\forall \pi$ . Therefore, we can  
 129 regard the above loss function that maximizes the performance gap (Eq. 2) as an instantiation of the  
 130 ranking-loss. We will refer to the implicit ranking between agent and the expert  $\rho^E \succeq \rho^\pi$  as vanilla  
 131 rankings and this variant of the ranking-loss function as the *supremum-loss*.

132 **Stackelberg Games:** A Stackelberg game is a general-sum game between two agents where one agent  
 133 is set to be the leader and the other a follower. The leader in this game optimizes its objective under the  
 134 assumption that the follower will choose the best response for its own optimization objective. More  
 135 concretely, assume there are two players  $A$  and  $B$  with parameters  $\theta_A, \theta_B$  and corresponding losses  
 136  $\mathcal{L}_A(\theta_A, \theta_B)$  and  $\mathcal{L}_B(\theta_A, \theta_B)$ . A Stackelberg game solves the following bi-level optimization when  
 137  $A$  is the leader and  $B$  is the follower:  $\min_{\theta_A} \mathcal{L}_A(\theta_A, \theta_B^*(\theta_A))$  s.t.  $\theta_B^*(\theta_A) = \arg \min_{\theta} \mathcal{L}_B(\theta_A, \theta)$ .  
 138 [51] showed that casting model-based RL as an approximate Stackelberg game [6] leads to  
 139 performance benefits and reduces training instability in comparison to the commonly used GDA [56]  
 140 and Best Reponse (BR) [12] methods. [17, 69] prove convergence of Stackelberg games under  
 141 smooth player cost functions and show that they reduce the cycling behavior to find an equilibrium  
 142 and allow for better convergence.

## 143 4 A Ranking Game for Imitation Learning

144 In this section, we first formalize the notion of the proposed two-player general-sum ranking game  
 145 for imitation learning. We then propose a practical instantiation of the ranking game through a  
 146 novel ranking-loss ( $L_k$ ). The proposed ranking game gives us the flexibility to incorporate additional  
 147 rankings—both auto-generated (a form of data augmentation mentioned as ‘auto’ in Fig. 1) and  
 148 offline (‘pref’ in Fig. 1)—which improves learning efficiency. Finally, we discuss the Stackelberg  
 149 formulation for the two-player ranking game and discuss two algorithms that naturally arise depending  
 150 on which player is designated as the leader.

### 151 4.1 The Two-Player Ranking Game Formulation

152 We present a new framework, `rank-game`, for imitation learning which casts it as a general-sum  
 153 *ranking game* between two players — a reward and a policy.

$$\underbrace{\arg \max_{\pi \in \Pi} J(R; \pi)}_{\text{Policy Agent}} \quad \underbrace{\arg \min_{R \in \mathcal{R}} L(\mathcal{D}^p; R)}_{\text{Reward Agent}}$$

154 In this formulation, the policy agent maximizes the reward by interacting with the environment, and  
 155 the reward agent attempts to find a reward function that satisfies a set of pairwise behavior rankings  
 156 in the given dataset  $\mathcal{D}^p$ ; a reward function satisfies these rankings if  $\mathbb{E}_{\rho^{\pi^i}}[R(s,a)] \leq \mathbb{E}_{\rho^{\pi^j}}[R(s,a)]$ ,  
 157  $\forall \rho^{\pi^i} \preceq \rho^{\pi^j} \in \mathcal{D}^p$ , where  $\rho^{\pi^i}, \rho^{\pi^j}$  can be state-action or state visitations.

158 The dataset of pairwise behavior rankings  $\mathcal{D}^p$  can be comprised of the implicit ‘vanilla’ rankings  
 159 between the learning agent and the expert’s policy behaviors ( $\rho^\pi \preceq \rho^E$ ), giving us the classical  
 160 IRL methods when a specific ranking loss function – *supremum-loss* is used [28, 22, 34]. If  
 161 rankings are provided between trajectories, they can be reduced to the equivalent ranking between the  
 162 corresponding state-action/state visitations. In the case when  $\mathcal{D}^p$  comprises purely of offline trajectory  
 163 performance rankings then, under a specific ranking loss function (*Luce-shepard*), the ranking game

---

**Algorithm 1** Meta algorithm: rank-game (vanilla) for imitation
 

---

- 1: Initialize policy  $\pi_\theta^0$ , reward function  $R_\phi$ , empty dataset  $\mathcal{D}^\pi$ . empirical expert data  $\hat{\rho}^E$
  - 2: **for**  $t = 1..T$  iterations **do**
  - 3:   Collect empirical visitation data  $\hat{\rho}^{\pi^t}$  with  $\pi^t$  in the environment. Set  $\mathcal{D}^\pi = \{(\hat{\rho}^\pi \preceq \hat{\rho}^E)\}$
  - 4:   Train reward  $R_\phi$  to satisfy rankings in  $\mathcal{D}^\pi$  using ranking loss  $L_k$  in equation 3.
  - 5:   Optimize policy under the reward function:  $\pi_\theta^{t+1} \leftarrow \operatorname{argmax}_{\pi'} J(R_\phi; \pi')$
  - 6: **end for**
- 

164 reduces to prior reward inference methods like T-REX [10, 11, 9, 13]. Thus, the ranking game affords  
 165 us a broader perspective of imitation learning, going beyond only using expert demonstrations.

## 166 4.2 Ranking Loss $L_k$ for the Reward Agent

167 We use a *ranking-loss* to train the reward function—an objective that minimizes the distortion [31]  
 168 between the ground truth ranking for a pair of entities  $\{x, y\}$  and rankings induced by a parameterized  
 169 function  $R : \mathcal{X} \rightarrow \mathbb{R}$  for a pair of scalars  $\{R(x), R(y)\}$ . One type of such a ranking-loss is the  
 170 *supremum-loss* in the classical imitation learning setup.

171 We propose a class of ranking-loss functions  $L_k$  that attempt to induce a performance gap of  $k$  for all  
 172 behavior preferences in the dataset. Formally, this can be implemented with the regression loss:

$$L_k(\mathcal{D}^p; R) = \mathbb{E}_{(\rho^{\pi^i}, \rho^{\pi^j}) \sim \mathcal{D}^p} \left[ \mathbb{E}_{s, a \sim \rho^{\pi^i}} [(R(s, a) - 0)^2] + \mathbb{E}_{s, a \sim \rho^{\pi^j}} [(R(s, a) - k)^2] \right]. \quad (3)$$

173 where  $\mathcal{D}^p$  contains behavior pairs  $(\rho^{\pi^i}, \rho^{\pi^j})$  s.t.  $\rho^{\pi^i} \preceq \rho^{\pi^j}$ .

174 The proposed ranking loss allows for learning *bounded rewards with user-defined scale  $k$*  in the agent  
 175 and the expert visitations as opposed to prior works in Adversarial Imitation Learning [28, 20, 22].  
 176 Reward scaling has been known to improve learning efficiency in deep RL; a large reward scale can  
 177 make the optimization landscape less smooth [27, 24] and a small scale might make the action-gap  
 178 small and increase susceptibility to extrapolation errors [7]. In contrast to the *supremum* loss,  $L_k$   
 179 can also naturally incorporate rankings provided by additional sources by learning a reward function  
 180 satisfying all specified pairwise preferences. The following theorem characterizes the equilibrium of  
 181 the rank-game for imitation learning when  $L_k$  is used as the ranking-loss.

182 **Theorem 4.1.** (*Performance of the rank-game equilibrium pair*) Consider an equilibrium of the  
 183 imitation rank-game  $(\hat{\pi}, \hat{R})$ , such that the ranking loss  $L_k$  generalization error is bounded by  
 184  $2R_{max}^2 \epsilon_r$  and the policy is near-optimal with  $J(\hat{R}; \hat{\pi}) \geq J(\hat{R}; \pi) - \epsilon_\pi \forall \pi$ , then at this equilibrium  
 185 pair under the expert’s unknown reward function  $R_{gt}$  bounded in  $[0, R_{max}^E]$ :

$$|J(R_{gt}, \pi^E) - J(R_{gt}, \hat{\pi})| \leq \frac{4R_{max}^E \sqrt{\frac{(1-\gamma)\epsilon_\pi + 4R_{max} \sqrt{\epsilon_r}}{k}}}{1 - \gamma} \quad (4)$$

186 If reward is a state-only function and only expert observations are available, the same bound applies  
 187 to the LfO setting.

188 *Proof.* We defer the proof to Appendix A. □

189 **Theoretical properties:** We now discuss some theoretical  
 190 properties of  $L_k$ . Theorem 1 shows that rank-game has  
 191 an equilibrium with bounded performance gap with the  
 192 expert. An optimization step by the policy player, under a  
 193 reward function optimized by the reward player, is equiva-  
 194 lent to minimizing an  $f$ -divergence with the expert. Equiva-  
 195 lently, at iteration  $t$  in Algorithm 1:  $\max_{\pi^t} \mathbb{E}_{\rho^{\pi^t}} [R_t^*] -$   
 196  $\mathbb{E}_{\rho^{\pi^E}} [R_t^*] = \min_{\pi^t} D_f(\rho^{\pi^t} \parallel \rho^{\pi^E})$ . We elaborate on the  
 197 regret of this idealized algorithm in Appendix A. Theorem  
 198 1 suggests that large values of  $k$  can guarantee the agent’s  
 199 performance is close to the expert. In practice, we observe  
 200 intermediate values of  $k$  also preserve imitation equilib-  
 201 rium optimality with a benefit of promoting sample efficient learning (as an effect of reward scaling

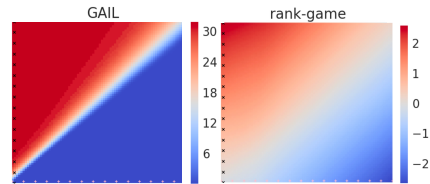


Figure 2: Figure shows learned reward function when agent and expert has a visitation shown by pink and black markers respectively. rank-game (auto) results in smooth reward functions more amenable to gradient-based policy optimization compared to GAIL.

described earlier). We discuss this observation further in Appendix D.9. `rank-game` naturally extends to the LfO regime under a state-only reward function where Theorem 4.1 results in a divergence bound between state-visitations of the expert and the agent. A state-only reward function is also a sufficient and necessary condition to ensure that we learn a dynamics-disentangled reward function [20].

$L_k$  can incorporate additional preferences that can help learn a regularized/shaped reward function that provides better guidance for policy optimization, reducing the exploration burden and increasing sample efficiency for IRL. A better-guided policy optimization is also expected to incur a lower  $\epsilon_\pi$ . However, augmenting the ranking dataset can lead to decrease in the intended performance gap ( $k_{eff} < k$ ) between the agent and the expert (Appendix A). This can loosen the bound in Eq 4 and lead to non-optimal imitation learning. We hypothesize that given informative preferences, decreased  $\epsilon_\pi$  can compensate potentially decreased intended performance gap  $k_{eff}$  to ensure near optimal imitation. In our experiments, we observe this hypothesis holds true; we enjoy sample efficiency benefits without losing any asymptotic performance. To leverage these benefits, we present two methods for augmenting the ranking dataset below and defer the implementation details to Appendix B.

#### 4.2.1 Generating the Ranking Dataset

**Reward loss w/ automatically generated rankings (auto):** In this method, we assume access to the behavior generating trajectories in the ranking dataset. For each pairwise comparison  $\rho_i \preceq \rho_j$  present in the dataset,  $L_k$  sets the regression targets for states in  $\rho_i$  to be 0 and for states visited by  $\rho_j$  to be  $k$ . Equivalently, we can rewrite minimizing  $L_k$  as regressing an input of trajectory  $\tau_i$  to vector  $\mathbf{0}$ , and  $\tau_j$  to vector  $k\mathbf{1}$  where  $\tau_i, \tau_j$  are trajectories that generate the behavior  $\rho_i, \rho_j$  respectively. We use the comparison  $\rho_i \preceq \rho_j$  to generate additional behavior rankings  $\rho_i \preceq \rho_{\lambda_1, ij} \preceq \rho_{\lambda_2, ij} \dots \preceq \rho_{\lambda_p, ij} \preceq \rho_j$  where  $0 < \lambda_1 < \lambda_2 < \dots < \lambda_p < 1$ . The behavior  $\rho_{\lambda_p, ij}$  is obtained by independently sampling the trajectories that generate the behaviors  $\rho_i, \rho_j$  and taking convex combinations i.e  $\tau_{\lambda_p, ij} = \lambda_p \tau_i + (1 - \lambda_p) \tau_j$  and their corresponding reward regressions targets are given by  $\lambda_p \mathbf{0} + (1 - \lambda_p) k \mathbf{1}$ .

This form of data augmentation can be interpreted as mixup [68] regularization in the trajectory space. Mixup has been shown to improve generalization and adversarial robustness [25, 68] by regularizing the first and second order gradients of the parameterized function. Following the general principle of using a smoothed objective with respect to inputs to obtain effective gradient signals, explicit smoothing in the trajectory-space can also help reduce the policy optimization error  $\epsilon_\pi$ . A didactic example showing rewards learned using this method is shown in Figure 2. In a special case when the expert’s unknown reward function is linear in observations, these rankings reflect the true underlying rankings of behaviors.

**Reward loss w/ offline annotated rankings (pref):** Another way of increasing learning efficiency is augmenting the ranking dataset containing the vanilla ranking ( $\rho^\pi \preceq \rho^E$ ) with offline annotated rankings. These rankings may be provided by a human observer or obtained using an offline dataset of behaviors with annotated reward information, similar to the datasets used in offline RL [19, 41]. We combine offline rankings by using a weighted loss between  $L_k$  for satisfying vanilla rankings ( $\rho^\pi \preceq \rho^E$ ) and offline rankings, grounded by an expert. Providing offline rankings alone that are sufficient to explain the reward function of the expert [10] is often a difficult task and the number of offline preferences required depends on the complexity of the environment. In the LfO setting, learning from an expert’s state visitation alone can be a hard problem due to exploration requirements [36]. This ranking-loss combines the benefits of using preferences to shape the reward function and guide policy improvement while using the expert to guarantee near-optimal performance.

#### 4.3 Optimizing the Two-Player General-Sum Ranking Game as a Stackelberg Game

Solving the ranking-game in the Stackelberg setup allows us to propose two different algorithms depending on which agent is set to be the leader and utilize the learning stability and efficiency afforded by the formulation as studied in [51, 69, 17].

**Policy as leader (PAL):** Choosing policy as the leader implies the following optimization:

$$\max_{\pi} \left\{ J(\hat{R}; \pi) \text{ s.t. } \hat{R} = \arg \min_R L(\mathcal{D}^\pi; R) \right\} \quad (5)$$

250 **Reward as leader (RAL):** Choosing reward as the leader implies the following optimization:

$$\min_{\hat{R}} \left\{ L(\mathcal{D}^\pi; \hat{R}) \text{ s.t. } \pi = \arg \max_{\pi} J(\hat{R}; \pi) \right\} \quad (6)$$

251 We follow the first order gradient approximation for leader’s update from previous work [51] to de-  
 252 velop practical algorithms. This strategy has been proven to be effective and avoids the computational  
 253 complexity of calculating the implicit Jacobian term ( $d\theta_B^*/d\theta_A$ ). PAL updates the reward to near con-  
 254 vergence on dataset  $\mathcal{D}^\pi$  ( $\mathcal{D}^\pi$  contains rankings generated using the current policy agent only  $\pi \preceq \pi^E$ )  
 255 and takes a few policy steps. Note that even after the first-order approximation, this optimization  
 256 strategy differs from GDA as often only a few iterations are used for training the reward even in hyper-  
 257 parameter studies like [46]. RAL updates the reward conservatively. This is achieved through aggregat-  
 258 ing the dataset of implicit rankings from all previous policies obtained during training. PAL’s strategy  
 259 of using on-policy data  $\mathcal{D}^\pi$  for reward training resembles that of methods including GAIL [28, 62],  
 260  $f$ -MAX [22], and  $f$ -IRL [45]. RAL uses the entire history of agent visitation to update the reward  
 261 function and resembles methods such as apprenticeship learning and DAC [1, 37]. PAL and RAL  
 262 bring together two seemingly different algorithm classes under a unified Stackelberg game viewpoint.

## 263 5 Experimental Results

264 We compare `rank-game` against state-of-the-art LfO and LfD approaches on MuJoCo benchmarks  
 265 having continuous state and action spaces. The LfO setting is more challenging since no actions are  
 266 available, and is a crucial imitation learning problem that can be used in cases where action modalities  
 267 differ between the expert and the agent, such as in robot learning. We focus on the LfO setting in this  
 268 section and defer the LfD experiments to Appendix D.2. We denote the imitation learning algorithms  
 269 that use the proposed ranking-loss  $L_k$  from Section 4.2 as RANK-`{PAL, RAL}`. We refer to the  
 270 `rank-game` variants which use automatically generated rankings and offline preferences as (auto)  
 271 and (pref) respectively following Section 4.2. In all our methods, we rely on an off-policy model-free  
 272 algorithm, Soft Actor-Critic (SAC) [26], for updating the policy agent.

273 We design experiments to answer the following questions:

- 274 1. *Asymptotic Performance and Sample Efficiency:* Is our method able to achieve near-expert  
 275 performance given a limited number (1) of expert observations? Can our method learn using fewer  
 276 environment interactions than prior state-of-the-art imitation learning (LfO) methods?
- 277 2. *Utility of preferences for imitation learning:* Current LfO methods struggle to solve a number  
 278 of complex manipulation tasks with sparse success signals. Can we leverage offline annotated  
 279 preferences through `rank-game` in such environments to achieve near-expert performance?
- 280 3. *Choosing between PAL and RAL methods:* Can we characterize the benefits and pitfalls of each  
 281 method, and determine when one method is preferable over the other?
- 282 4. *Ablations for the method components:* Can we establish the importance of hyperparameters and  
 283 design decisions in our experiments?

284 **Baselines:** We compare RANK-PAL and RANK-RAL against 6 representative LfO approaches that  
 285 covers a spectrum of on-policy and off-policy model-free methods from prior work: GAIfo [62, 28],  
 286 DACfo [37], BCO [61],  $f$ -IRL [45] and recently proposed OPOLO [71] and IQLearn [21]. We do  
 287 not assume access to expert actions in this setting. Our LfD experiments compare to the IQLearn [21],  
 288 DAC [37] and BC baselines. Detailed description for baselines can be found in Appendix D.2.

### 289 5.1 Asymptotic Performance and Sample Efficiency

290 In this section, we compare RANK-PAL(auto) and RANK-RAL(auto) to baselines on a set of MuJoCo  
 291 locomotion tasks of varying complexities: `Swimmer-v2`, `Hopper-v2`, `HalfCheetah-v2`,  
 292 `Walker2d-v2`, `Ant-v2` and `Humanoid-v2`. In this experiment, we provide one expert trajec-  
 293 tory for all methods and do not assume access to any offline annotated rankings.

294 **Asymptotic Performance:** Table 2 shows that both `rank-game` methods are able to reach near-  
 295 expert asymptotic performance with a single expert trajectory. BCO shows poor performance which

Env	Hopper	HalfCheetah	Walker	Ant	Humanoid
BCO	20.10±2.15	5.12±3.82	4.00±1.25	12.80±1.26	3.90±1.24
GAIFo	81.13± 9.99	13.54±7.24	83.83±2.55	20.10±24.41	3.93±1.81
DACfO	94.73±3.63	85.03±5.09	54.70±44.64	86.45±1.67	19.31±32.19
$f$ -IRL	97.45± 0.61	96.06±4.63	<b>101.16±1.25</b>	71.18±19.80	77.93±6.372
OPOLO	89.56±5.46	88.92±3.20	79.19±24.35	93.37± 3.78	24.87±17.04
RANK-PAL(ours)	87.14± 16.14	94.05±3.59	93.88±0.72	<b>98.93±1.83</b>	<b>96.84±3.28</b>
RANK-RAL(ours)	<b>99.34±0.20</b>	<b>101.14±7.45</b>	93.24±1.25	93.21±2.98	94.45±4.13
Expert	100.00± 0	100.00± 0	100.00± 0	100.00± 0	100.00± 0
( $ S $ , $ A $ )	(11, 3)	(17, 6)	(17, 6)	(111, 8)	(376, 17)

Table 2: Asymptotic normalized performance of LfO methods at 2 million timesteps on MuJoCo locomotion tasks. The standard deviation is calculated with 5 different runs [each averaging over 10 trajectory returns](#). For unnormalized score and more details, check Appendix D. We omit IQlearn due to poor performance.

296 can be attributed to the compounding error problem arising from its behavior cloning strategy. GAIFo  
 297 and DACfO use GDA for optimization with a supremum loss and show high variance in their  
 298 asymptotic performance whereas rank-game methods are more stable and low-variance.

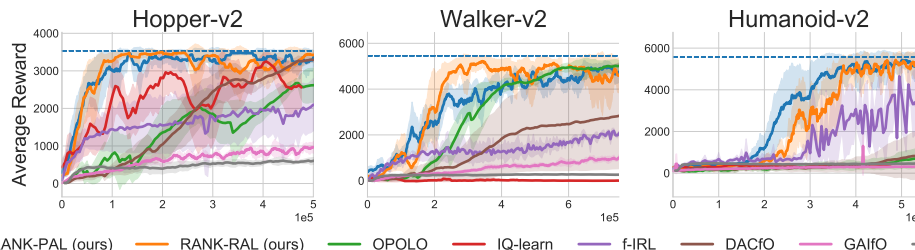


Figure 3: Comparison of performance on OpenAI gym benchmark tasks. The shaded region represents standard deviation across 5 random runs. RANK-PAL and RANK-RAL substantially outperform the baselines in sample efficiency. Complete set of results can be found in Appendix D.1

299 **Sample Efficiency:** Figure 3 shows that RANK-RAL and RANK-PAL are among the most sample  
 300 efficient methods for the LfO setting, outperforming the recent state-of-the-art method OPOLO [71]  
 301 by a significant margin. We notice that IQLearn fails to learn in the LfO setting. This experiment  
 302 demonstrates the benefit of the combined improvements of the proposed ranking-loss with automati-  
 303 cally generated rankings. Our method is also simpler to implement than OPOLO, as we require  
 304 fewer lines of code changes on top of SAC and need to maintain fewer parameterized networks  
 305 compared to OPOLO which requires an additional inverse action model to regularize learning.  
 306

## 307 5.2 Utility of Preferences in Imitation

308 Our experiments on complex manipulation environments—door opening with a parallel-  
 309 jaw gripper [70] and pen manipulation with a dexterous adroit hand [50] – reveal that none  
 310 of the prior LfO methods are able to imitate the expert even under increasing amounts of  
 311 expert data. This failure of LfO methods can be potentially attributed to the exploration  
 312 requirements of LfO compared to LfD [36], coupled with the sparse successes encountered  
 313 in these tasks, leading to poorly guided policy gradients. In these experiments, we show that  
 314 rank-game can incorporate additional information in the form of offline annotated rankings to  
 315 guide the agent in solving such tasks. These offline rankings are obtained by uniformly sampling  
 316 a small set of trajectories (10) from the replay buffer of SAC [26] labeled with a ground truth  
 317 reward function. We use a weighted ranking loss (pref) from Section 4.2.

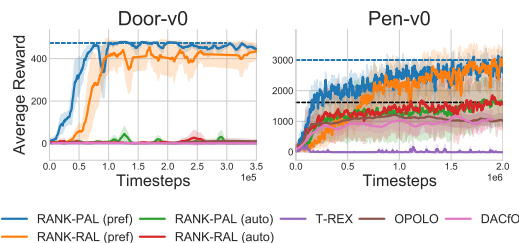


Figure 4: Offline annotated preferences can help solve LfO tasks in the complex manipulation environments Pen-v0 and Door, whereas prior LfO methods fail. Black dotted line shows asymptotic performance of RANK-PAL (auto) method.

323 Figure 4 shows that RANK-PAL/RAL(pref) method leveraging offline ranking is the only method that  
 324 can solve these tasks, whereas prior LfO methods and RANK-PAL/RAL(auto) with automatically



325 generated rankings struggle even after a large amount of training. We also point out that T-REX, a  
 326 method that learns using the preferences alone is unable to achieve near-expert performance, thereby  
 327 highlighting the benefits of learning from expert demonstrations alongside a set of offline preferences.

### 328 5.3 Comparing PAL and RAL

329 PAL uses the agent’s current visitation for re-  
 330 ward learning, whereas RAL learns a reward  
 331 consistent with all rankings arising from the history  
 332 of the agent’s visitation. These properties  
 333 can present certain benefits depending on the  
 334 task setting. To test the potential benefits of PAL  
 335 and RAL, we consider two non-stationary imitation  
 336 learning problems, similar to [50] – one in  
 337 which the expert changes its intent and the other  
 338 where dynamics of the environment change dur-  
 339 ing training in the Hopper-v2 locomotion task.  
 340 For changing intent, we present a new set of  
 341 demonstrations where the hopper agent hops  
 342 backwards rather than forward. For changing  
 343 environment dynamics, we increase the mass of  
 344 the hopper agent by a factor of 1.2. Changes are  
 345 introduced at  $1e5$  time steps during training at which point we notice a sudden performance drop.  
 346 In Figure 5 (left), we notice that PAL adapts faster to intent changes, whereas RAL needs to unlearn  
 347 the rankings obtained from the agent’s history and takes longer to adapt. Figure 5 (right) shows that  
 348 RAL adapts faster to the changing dynamics of the system, as it has already learned a good global  
 349 notion of the dynamics-disentangled reward function in the LfO setting, whereas PAL only has a local  
 350 understanding of reward as a result of using ranking obtained only from the agent’s current visitation.

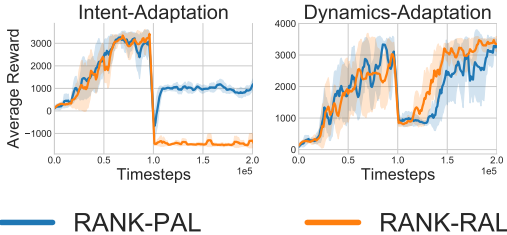


Figure 5: We compare the relative strengths of PAL and RAL. Left plot shows a comparison when the goal is changed, and right plot shows a comparison when dynamics of the environment is changed. These changes occur at  $1e5$  timesteps into training. PAL adapts faster to changing intent and RAL adapts faster to changing dynamics.

351 **Ablation of Method Components:** Appendix D contains eight additional experiments to study the  
 352 importance of hyperparameters and design decisions. Our ablations validate the importance of using  
 353 automatically generated rankings, the benefit of ranking loss over *supremum* loss, and sensitivity to  
 354 hyperparameters like the intended performance gap  $k$ , policy iterations, and the reward regularizer.

## 355 6 Conclusion

356 In this work, we present a new framework for imitation learning that treats imitation as a two-player  
 357 ranking-game between a policy and a reward function. Unlike prior works in imitation learning, the  
 358 ranking game allows incorporation of rankings over suboptimal behaviors to aid policy learning. We  
 359 instantiate the ranking game by proposing a novel ranking loss which guarantees agent’s performance  
 360 to be close to expert for imitation learning. Our experiments on simulated MuJoCo tasks reveal that  
 361 utilizing additional ranking through our proposed ranking loss leads to improved sample efficiency  
 362 for imitation learning, outperforming prior methods by a significant margin and solving some tasks  
 363 which were unsolvable by previous LfO methods.

364 **Limitations and Negative Societal Impacts:** Preferences obtained in real world are usually  
 365 noisy [40, 32, 8] and one limitation of *rank-game* is that it does not suggest a way to handle noisy  
 366 preferences. Second, *rank-game* proposes modifications to learn a reward function amenable to  
 367 policy optimization but these hyperparameters are set manually. Future work can explore methods to  
 368 automate learning such reward functions. Third, despite learning effective policies we observed that  
 369 we do not learn reusable robust reward functions [45]. Negative Societal Impact: Imitation learning  
 370 can cause harm if given demonstrations of harmful behaviors, either accidentally or purposefully.  
 371 Furthermore, even when given high-quality demonstrations of desirable behaviors, our algorithm does  
 372 not provide guarantees of performance, and thus could cause harm if used in high-stakes domains  
 373 without sufficient safety checks on learned behaviors.

374 **References**

- 375 [1] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning.  
376 In *Proceedings of the twenty-first international conference on Machine learning*, page 1, 2004.
- 377 [2] Joshua Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- 378 [3] Riad Akrouf, Marc Schoenauer, and Michele Sebag. Preference-based policy learning. In *Joint*  
379 *European Conference on Machine Learning and Knowledge Discovery in Databases*, pages  
380 12–27. Springer, 2011.
- 381 [4] Syed Mumtaz Ali and Samuel D Silvey. A general class of coefficients of divergence of one  
382 distribution from another. *Journal of the Royal Statistical Society: Series B (Methodological)*,  
383 28(1):131–142, 1966.
- 384 [5] Oleg Arenz and Gerhard Neumann. Non-adversarial imitation learning and its connections to  
385 adversarial methods. *arXiv preprint arXiv:2008.03525*, 2020.
- 386 [6] Tamer Başar and Geert Jan Olsder. *Dynamic noncooperative game theory*. SIAM, 1998.
- 387 [7] Marc G Bellemare, Georg Ostrovski, Arthur Guez, Philip Thomas, and Rémi Munos. Increasing  
388 the action gap: New operators for reinforcement learning. In *Proceedings of the AAAI*  
389 *Conference on Artificial Intelligence*, volume 30, 2016.
- 390 [8] Erdem Bıyık, Dylan P Losey, Malayandi Palan, Nicholas C Landolfi, Gleb Shevchuk, and  
391 Dorsa Sadigh. Learning reward functions from diverse sources of human feedback: Optimally  
392 integrating demonstrations and preferences. *The International Journal of Robotics Research*,  
393 page 02783649211041652, 2021.
- 394 [9] Daniel Brown, Russell Coleman, Ravi Srinivasan, and Scott Niekum. Safe imitation learning  
395 via fast bayesian reward inference from preferences. In *International Conference on Machine*  
396 *Learning*, pages 1165–1177. PMLR, 2020.
- 397 [10] Daniel S. Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating be-  
398 yond suboptimal demonstrations via inverse reinforcement learning from observations. *ArXiv*,  
399 abs/1904.06387, 2019.
- 400 [11] Daniel S Brown, Wonjoon Goo, and Scott Niekum. Better-than-demonstrator imitation learning  
401 via automatically-ranked demonstrations. In *Conference on robot learning*, pages 330–359.  
402 PMLR, 2020.
- 403 [12] Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university  
404 press, 2006.
- 405 [13] Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration  
406 via self-supervised reward regression. *arXiv preprint arXiv:2010.11723*, 2020.
- 407 [14] Paul Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei.  
408 Deep reinforcement learning from human preferences. *arXiv preprint arXiv:1706.03741*, 2017.
- 409 [15] Imre Csiszár. Information-type measures of difference of probability distributions and indirect  
410 observation. *studia scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- 411 [16] Ashley Edwards, Himanshu Sahni, Yannick Schroecker, and Charles Isbell. Imitating latent  
412 policies from observation. In *International Conference on Machine Learning*, pages 1755–1763.  
413 PMLR, 2019.
- 414 [17] Tanner Fiez, Benjamin Chasnov, and Lillian J Ratliff. Convergence of learning dynamics in  
415 stackelberg games. *arXiv preprint arXiv:1906.01217*, 2019.

- 416 [18] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal  
417 control via policy optimization. In *International conference on machine learning*, pages 49–58.  
418 PMLR, 2016.
- 419 [19] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
420 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 421 [20] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse  
422 reinforcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
- 423 [21] Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn:  
424 Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34,  
425 2021.
- 426 [22] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence mini-  
427 mization perspective on imitation learning methods. In *Conference on Robot Learning*, pages  
428 1259–1277. PMLR, 2020.
- 429 [23] Gustavo Gilardoni. On the minimum f-divergence for given total variation. *Comptes Rendus*  
430 *Mathematique - C R MATH*, 343:763–766, 12 2006.
- 431 [24] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedfor-  
432 ward neural networks. In *Proceedings of the thirteenth international conference on artificial*  
433 *intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 434 [25] Hongyu Guo, Yongyi Mao, and Richong Zhang. Mixup as locally linear out-of-manifold  
435 regularization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33,  
436 pages 3714–3722, 2019.
- 437 [26] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
438 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
439 *conference on machine learning*, pages 1861–1870. PMLR, 2018.
- 440 [27] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David  
441 Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on*  
442 *artificial intelligence*, volume 32, 2018.
- 443 [28] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*  
444 *information processing systems*, 29:4565–4573, 2016.
- 445 [29] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. In *The*  
446 *collected works of Wassily Hoeffding*, pages 409–426. Springer, 1994.
- 447 [30] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward  
448 learning from human preferences and demonstrations in atari. *arXiv preprint arXiv:1811.06521*,  
449 2018.
- 450 [31] Rishabh Iyer and Jeff Bilmes. The submodular bregman and lovász-bregman divergences with  
451 applications: Extended version. Citeseer, 2012.
- 452 [32] Hong Jun Jeon, Smitha Milli, and Anca Dragan. Reward-rational (implicit) choice: A unifying  
453 formalism for reward learning. *Advances in Neural Information Processing Systems*, 33:4415–  
454 4426, 2020.
- 455 [33] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning.  
456 In *In Proc. 19th International Conference on Machine Learning*. Citeseer, 2002.
- 457 [34] Liyiming Ke, Sanjiban Choudhury, Matt Barnes, Wen Sun, Gilwoo Lee, and Siddhartha Srimi-  
458 vasa. Imitation learning as f-divergence minimization. In *Algorithmic Foundations of Robotics*  
459 *XIV: Proceedings of the Fourteenth Workshop on the Algorithmic Foundations of Robotics 14*,  
460 pages 313–329. Springer International Publishing, 2021.

- 461 [35] Michael Kearns and Satinder Singh. Finite-sample convergence rates for q-learning and indirect  
462 algorithms. *Advances in neural information processing systems*, 11, 1998.
- 463 [36] Rahul Kidambi, Jonathan Chang, and Wen Sun. Mobile: Model-based imitation learning from  
464 observation alone. *Advances in Neural Information Processing Systems*, 34, 2021.
- 465 [37] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan  
466 Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in  
467 adversarial imitation learning. *arXiv preprint arXiv:1809.02925*, 2018.
- 468 [38] Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribu-  
469 tion matching. *arXiv preprint arXiv:1912.05032*, 2019.
- 470 [39] Victoria Krakovna. Specification gaming examples in ai. *Available at vkrakovna.wordpress.*  
471 *com*, 2018.
- 472 [40] Minae Kwon, Erdem Biyik, Aditi Talati, Karan Bhasin, Dylan P Losey, and Dorsa Sadigh.  
473 When humans aren't optimal: Robots that collaborate with risk-aware humans. In *2020 15th*  
474 *ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 43–52. IEEE,  
475 2020.
- 476 [41] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
477 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 478 [42] Friedrich Liese and Igor Vajda. On divergences and informations in statistics and information  
479 theory. *IEEE Transactions on Information Theory*, 52(10):4394–4412, 2006.
- 480 [43] Fangchen Liu, Zhan Ling, Tongzhou Mu, and Hao Su. State alignment-based imitation learning.  
481 *arXiv preprint arXiv:1911.10947*, 2019.
- 482 [44] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*,  
483 volume 1, page 2, 2000.
- 484 [45] Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Benjamin Eysenbach. f-irl:  
485 Inverse reinforcement learning via state marginal matching. *arXiv preprint arXiv:2011.04709*,  
486 2020.
- 487 [46] Manu Orsini, Anton Raichuk, Léonard Hussenot, Damien Vincent, Robert Dadashi, Sertan  
488 Girgin, Matthieu Geist, Olivier Bachem, Olivier Pietquin, and Marcin Andrychowicz. What  
489 matters for adversarial imitation learning? *Advances in Neural Information Processing Systems*,  
490 34, 2021.
- 491 [47] Malayandi Palan, Nicholas C Landolfi, Gleb Shevchuk, and Dorsa Sadigh. Learning reward func-  
492 tions by integrating human demonstrations and preferences. *arXiv preprint arXiv:1906.08928*,  
493 2019.
- 494 [48] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation.  
495 *Neural computation*, 3(1):88–97, 1991.
- 496 [49] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*.  
497 John Wiley & Sons, 2014.
- 498 [50] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel  
499 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement  
500 learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- 501 [51] Aravind Rajeswaran, Igor Mordatch, and Vikash Kumar. A game theoretic framework for model  
502 based reinforcement learning. In *ICML*, 2020.

- 503 [52] Siddharth Reddy, Anca D Dragan, and Sergey Levine. Sqil: Imitation learning via reinforcement  
504 learning with sparse rewards. *arXiv preprint arXiv:1905.11108*, 2019.
- 505 [53] Alfréd Rényi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley*  
506 *Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory*  
507 *of Statistics*, pages 547–561. University of California Press, 1961.
- 508 [54] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and  
509 structured prediction to no-regret online learning. In *Proceedings of the fourteenth interna-*  
510 *tional conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and  
511 Conference Proceedings, 2011.
- 512 [55] Dorsa Sadigh, Anca D Dragan, Shankar Sastry, and Sanjit A Seshia. Active preference-based  
513 learning of reward functions. 2017.
- 514 [56] Florian Schäfer and Anima Anandkumar. Competitive gradient descent. *arXiv preprint*  
515 *arXiv:1905.12103*, 2019.
- 516 [57] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*  
517 *algorithms*. Cambridge university press, 2014.
- 518 [58] Wen Sun, Anirudh Vemula, Byron Boots, and Drew Bagnell. Provably efficient imitation  
519 learning from observation alone. In *International conference on machine learning*, pages  
520 6036–6045. PMLR, 2019.
- 521 [59] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,  
522 2018.
- 523 [60] Gokul Swamy, Sanjiban Choudhury, J Andrew Bagnell, and Steven Wu. Of moments and match-  
524 ing: A game-theoretic framework for closing the imitation gap. In *International Conference on*  
525 *Machine Learning*, pages 10022–10032. PMLR, 2021.
- 526 [61] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv*  
527 *preprint arXiv:1805.01954*, 2018.
- 528 [62] Faraz Torabi, Garrett Warnell, and Peter Stone. Generative adversarial imitation from observa-  
529 tion. *arXiv preprint arXiv:1807.06158*, 2018.
- 530 [63] Faraz Torabi, Garrett Warnell, and Peter Stone. Recent advances in imitation learning from  
531 observation. *arXiv preprint arXiv:1905.13566*, 2019.
- 532 [64] Igor Vajda. Note on discrimination information and variation (corresp.). *IEEE Transactions on*  
533 *Information Theory*, 16(6):771–773, 1970.
- 534 [65] Aaron Wilson, Alan Fern, and Prasad Tadepalli. A bayesian approach for policy learning from  
535 trajectory preference queries. *Advances in neural information processing systems*, 25:1133–  
536 1141, 2012.
- 537 [66] Tian Xu, Ziniu Li, and Yang Yu. Error bounds of imitating policies and environments for  
538 reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- 539 [67] Chao Yang, Xiaojian Ma, Wenbing Huang, Fuchun Sun, Huaping Liu, Junzhou Huang, and  
540 Chuang Gan. Imitation learning from observations by minimizing inverse dynamics disagree-  
541 ment. *arXiv preprint arXiv:1910.04417*, 2019.
- 542 [68] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond  
543 empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 544 [69] Liyuan Zheng, Tanner Fiez, Zane Alumbaugh, Benjamin Chasnov, and Lillian J Ratliff.  
545 Stackelberg actor-critic: Game-theoretic reinforcement learning algorithms. *arXiv preprint*  
546 *arXiv:2109.12286*, 2021.

- 547 [70] Yuke Zhu, Josiah Wong, Ajay Mandlekar, and Roberto Martín-Martín. robosuite: A modular  
548 simulation framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*,  
549 2020.
- 550 [71] Zhuangdi Zhu, Kaixiang Lin, Bo Dai, and Jiayu Zhou. Off-policy imitation learning from  
551 observations. *Advances in Neural Information Processing Systems*, 33, 2020.
- 552 [72] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy  
553 inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.