# Swarm Inverse Reinforcement Learning for Biological Systems

1st Xin Yu
*School of Computer Science and Engineering*
*Beihang University*
Beijing, China
nlsdeyuxin@buaa.edu.cn

2nd Wenjun Wu
*Institute of Artificial Intelligence*
*Beihang University*
Beijing, China
wwj09315@buaa.edu.cn

3rd Pu Feng
*School of Computer Science and Engineering*
*Beihang University*
Beijing, China
fengpu@buaa.edu.cn

4th Yongkai Tian
*School of Computer Science and Engineering*
*Beihang University*
Beijing, China
tianyk@buaa.edu.cn

*Abstract*—**Complex global behavior can emerge from local interactions in biological systems. Many models have been introduced to describe the interaction rules of biological individuals. Nonetheless, most research efforts cannot capture the inner cognitive and sequential decision process of individual animals in their swarms. In this paper, we formulate this problem as homogeneous Markov game and focus on identifying the potential reward function of individual animals so as to understand their collective behaviors. We propose an inverse reinforcement learning method PS-AIRL specifically for biological systems, where the parameter sharing paradigm is combined with a deep inverse reinforcement learning. Theoretical analysis and experimental evaluation show that PS-AIRL can learn the policy and the reward function from collective behavior demonstrations. Moreover, our methods can be applied to a wide range of biological behavioral studies.**

*Index Terms*—**collective intelligence, imitation learning, multi-agent system, reinforcement learning**

## I. INTRODUCTION

In biological systems, many identical agents interact with each other to achieve a common survival goal and demonstrate collective behavior. Collective behavior can be found in many species such as flocking birds and collective swimming of microorganisms, which bring them advantages over individual behavior including increased diffusion, faster transport, and faster foraging. Various models have been proposed to understand how individual animal decisions based on their local observations can emerge as collective behaviors at the global level. A common conclusion from such research confirms that even simple rules can produce complex behaviors on a global

scale [1]. Thus, more investigations must be conducted to reveal the relationship between global behavior and agent-level interactions. Such a research topic is not only a biological interest but also key to cyber-physical system applications such as distributed sensor networks and multi robots systems.

There are two major approaches to study animal interactions in collective behaviors including the agent-based modeling (ABM) and machine learning. ABM is widely used to study complex collective dynamics that consist of autonomous agents interacting with each other [2], [3]. Vicsek model is a seminal agent-based model describing a flocking phase transition. However, constructing such a model needs to propose a set of candidate local rules and compare the results simulated by rules with the desired outcome until an adequate configuration is found. These individual rules are often designed based on domain knowledge and require careful adjustment for parameters [4]. While ABM provides an insightful view of collective behavior, Vicsek models based on the abstraction of particles for individuals cannot capture complex inner cognitive and biophysical factors that are necessary for accurately describing individual behavior [5]. To tackle these issues, machine learning techniques were proposed to model collective behaviors in a data-driven way diminishing the requirement of domain knowledge [6]. For example, deep neural networks were trained to predict the future turning side of a zebrafish and attained higher accuracy than previous agent-based models [7]. Though these models exhibit high accuracy for behavioral prediction, they haven't reveal latent cognitive decision mechanisms to drive individual animals to respond in specific ways towards various swarm behaviors [8].

A biological system consisting of multiple agents is mainly formulated as Markov game. Markov game assumes that each agent in the system follows a policy aiming to maximize

its internal reward. The reward function is a thorough characterization of the interacting tendency for the agents [9]. A promising way to understand collective interaction is to recover the potential reward function each agent follows. Inverse reinforcement learning (IRL) provides a data-driven approach to approximate the reward function from collective behavior data [10]. Recent work adopts IRL to find interaction rules in collectives [11], [12]. Unfortunately, these methods only consider the single-agent IRL which is not suitable for collective settings due to the non-stationary environment for individual agents [13]. The biological swarm system is composed of a large number of homogeneous agents and keeps high dimension observation space and action space. Therefore, reconstructing the reward function for biological systems needs to exploiting multi-agent IRL methods. Previous multi-agent IRL methods can't support the scale of biological swarms because of their high computational overhead and poor algorithmic convergence in practice [14], [15]. The remedy to this issue requires considerations for the homogeneity and locality of a biological swarm. Another research efforts in [16] attempt to extend inverse reinforcement learning to homogeneous multi-agent systems. But this algorithm assume simple representations of strategy and reward function, thus not applicable for capturing in high dimensional features in animals' sequential decision mechanisms.

In contrast to previous methods, we propose multi-agent inverse reinforcement learning methods specially adapted to biological systems. Identical individuals interact locally in a biological system contribute to the characteristic of homogeneity and locality. Considering the aforementioned characteristics, we introduce an IRL solution called parameter sharing adversarial inverse reinforcement learning (PS-AIRL) that can solve the high dimensional problem by combining the parameter sharing paradigm and the deep IRL methods. There are two objectives of PS-AIRL: (1) policy imitation, learning policies by imitating biological systems. (2) reward reconstruction, recovering reward functions that induce collective behaviors. We make theoretical analysis and experimental evaluation for the PS-AIRL. The contributions of this paper are summarized as follows:

- By considering homogeneity and locality of biological swarms, we propose a parameter sharing deep IRL method specially adapted to swarm systems.
- We present theoretical analysis for the parameter sharing mechanism in multi-agent reinforcement learning to illustrate the effectiveness of the PS-AIRL framework.
- We conduct a variety of experiments in two major swarm scenarios. Experimental results show that the PS-AIRL can explain and accurately reproduce the collective behavior of a swarm system.

## II. Methods

In this section, we first present the notation definition for the swarm system. Secondly, we propose the parameter-sharing learning paradigm with theoretical analysis. Finally, we describe our swarm inverse reinforcement learning method in detail.

### A. Problem formulation

A biological swarm system often exhibits the characteristics of homogeneity and locality.

- Homogeneity: All agents in the system carry a common architecture (i.e. The same observation space and action space)
- Locality: The agents can observe only parts of the system within a certain range. Their decisions depend on their current neighborhood only.

In principle, any system with these properties can be described as a Markov game [17]. Given the homogeneity property of swarm systems, we can further adopt the swarMDPs that explicitly implements a homogeneous Markov game architecture [16]. The SwarMDP framework is defined as a tuple $(N, S, O, A, R, T, \pi, \xi)$.

- $N$ is the number of agents in the system.
- $S, O, A$ are sets of local states, observations and actions respectively.
- $R : O \to R$ is an agent-level reward function.
- $T : S^N \times A^N \times S^N \to R$ is the global transition model of the system. The system reaches state $\tilde{s} = (\tilde{s}^{(1)}, \ldots, \tilde{s}^{(N)})$ when the agents perform the joint action $a = (a^{(1)}, \ldots, a^{(N)})$ at state as $T(\tilde{s}|s,a)$, where $s^{(n)}, \tilde{s}^{(n)} \in S$ and $a^{(n)} \in A_1$ represent the local states and the local action of agent $n$, respectively.
- $\pi : O \to A$ is the local policy.
- $\xi : S^N \to O^N$ is the observation model of the system.

The observation model $\xi$ indicates every agent's sensing capability, which defines their perception of a given system state $s \in S^N$. For example, in a flocking of birds, $\xi^{(n)}$ could represent a bird's local perception of its immediate neighbors.

### B. Parameter sharing for biological system

Actor critic algorithms are a class of model-free RL algorithms which can be used to train the policy network in the inverse reinforcement learning algorithm [18]. In a multi-agent partially observable setting, the simplest AC algorithm defines a policy loss in equation(1):

$$L(\phi_i) = -\log \pi_{\phi_i}(a_t^i|o_t^i)((r_t^i + \gamma V_{\theta_i}(o_{t+1}^i) - V_{\theta_i}(o_t^i)) \quad (1)$$

The value loss for agent i in the multi-agent partially observable setting can be defined in equation (2):

$$L(\theta_i) = \|V_{\theta_i}(o_t^i) - y_i\|^2 \text{ with } y_i = r_t^i + \gamma V_{\theta_i}(o_{t+1}^i) \quad (2)$$

The number of animals in a biological swarm system is often very large. To efficiently train policy networks for such a large homogeneous multi-agent swarm, it is common to adopt parameter sharing technique, where all agents share the same parameter in their policy networks [19]. Though widely used, the effectiveness of the parameter sharing paradigm has not been confirmed by theoretical analysis and is often thought of as an implementation detail [20].

From the perspective of one agent, other agents are perceived as part of the environment. In each episode, every agent interacts with the environment to generate its trajectories. Traditionally, an agent learns its policy network and value network from its own experience according to the loss functions specified in (1) and (2) respectively. To accelerate training speed, Christianos et al. [21] applies experience sharing by combining the gradients of different agents to update the policy separately. Their main theoretical results indicating that the policy network and the value network can be updated by their own trajectories as well as the experience of other agents in a multi-agent scenario. The policy network can be updated in equation (3):

$$L(\phi_i) = -\log \pi_{\phi_i}(a_t^i|o_t^i)(r_t^i + \gamma V_{\theta_i}(o_{t+1}^i) - V_{\theta_i}(o_t^i)) - \lambda \cdot$$
$$\sum_{k \neq i} \frac{\pi_{\phi_i}(a_t^k|o_t^k)}{\pi_{\phi_k}(a_t^k|o_t^k)} \log \pi_{\phi_i}(a_t^k|o_t^k)(r_t^k + \gamma V_{\theta_i}(o_{t+1}^k) - V_{\theta_i}(o_t^k))$$
(3)

The value network can be updated in equation (4).The hyper-parameter $\lambda$ weights the experience of other agents.

$$L(\theta_i) = \|V_{\theta_i}(o_t^i) - y_i^i\|^2 + \lambda \sum_{k \neq i} \frac{\pi_{\theta_i}(a_t^k|o_t^k)}{\pi_{\phi_k}(a_t^k|o_t^k)} \|V_{\theta_i}(o_t^k) - y_k^i\|^2$$
(4)

The intuition of parameter sharing is that all the agents share the same network, which is learned from all the trajectories. All the agents execute the same policy illustrated in equation(5):

$$\pi_{\phi_i}(a_t^k|o_t^k) = \pi_{\phi_k}(a_t^k|o_t^k)$$
(5)

The loss function (3) and (4) of the policy network and the value network are reduced to (6) and (7), which means that the policy network and the value network can be updated by all the trajectories generated from the biological system.

$$L(\phi_i) = -\sum_k \log \pi_{\phi_i}(a_t^k|o_t^k)(r_t^k + \gamma V_{\theta_i}(o_{t+1}^k) - V_{\theta_i}(o_t^k))$$
(6)

$$L(\theta_i) = \sum_k \|V_{\theta_i}(o_t^k) - y_k^i\|^2$$
(7)

During the reinforcement learning process in the environment, each agent executes the same policy network. Equation (6) and (7) shows that the sum of policy and value loss gradients are used to optimise the shared parameters.

*C. Swarm inverse reinforcement learning*

Adversarial Inverse Reinforcement Learning (AIRL) is a promising single-agent inverse reinforcement learning method based on the generative adversarial framework which consists of the generator and the discriminator [22] [23]. In AIRL, the agent interacts with the environment by executing a policy network $\pi$ to generate trajectories $\tau_j$. And these trajectories $\tau_j$ must be compared against the expert trajectories $\tau_E$ by the discriminator $D_\omega$ parameterised by $\omega$.

In swarm imitation learning settings, we do not have access to the reward function, but have demonstrations provided by experts (N expert agents in swarMDP). Fig. 1 shows the PS-AIRL framework, where we formulate biological swarm systems by swarMDP and extend the AIRL by sharing the generator and discriminator among all the agents. The goal of PS-AIRL is to infer the right reward function for the agents so as to mimic the experts' policies $\pi^E$.



Fig. 1. PS-AIRL framework for biological system.

The learning procedure for PS-AIRL is summarized in Algorithm 1. The demonstrations can be denoted as $\tau_E = \{\tau_j\}_{j=1}^M$, where $\tau_j = \{(s_j^t, \boldsymbol{a}_j^t)\}_{t=1}^T$ is the trajectory of animal $j$ collected from step 1 to step T. In the first step, we maintain a collective demonstrations $\tau_E$ collected by $a_t = \pi_E(a_t|s_t)$. PS-AIRL randomly initializes the policy network $\pi$ and the discriminator network $D_{\theta,\phi}$. Each agent interacts with the environment independently according to the current policy. The resulting trajectories denoted as $\vec{\tau}$ are sampled by different agents following the same policy $\pi_{\theta_k}$. The training process of discriminator $D_{\theta,\phi}$ in multi-agent IRL problem is the same as single-agent one. PS-AIRL trains $D_{\theta,\phi}$ via binary logistic regression to classify expert data $\tau_E$ from samples $\vec{\tau}$. After training the discriminator, it can update the reward function(Algorithm 1, Line 6). In each iteration, the algorithm trains the policy network $\pi$ by Algorithm 2.

---

**Algorithm 1** PS-AIRL

1: Input: expert trajectories $\tau_E \sim \pi_E$
2: initialize policy $\pi_{\theta_0}$ and discriminator $D_{\theta_0,\phi_0}$
3: **for** $k \leftarrow 1, 2, \ldots$ **do**
4:      Rollout trajectories for all agents $\vec{\tau} \sim \pi_{\theta_k}$
5:      Train $D_{\theta_k,\phi_k}$ via binary logistic regression to classify expert data $\tau_E$ from samples $\vec{\tau}$.
6:      Generating reward $r_{\theta_k,\phi_k}$

$$r_{\theta_k,\phi_k} \rightarrow \log D_{\theta_k,\phi_k} - \log(1 - D_{\theta_k,\phi_k})$$

7:      Updating policy $\pi_k$ with respect to $r_{\theta_k,\phi_k}$ by PS-PPO
8: **end for**

---

Proximal Policy Optimization (PPO) is an Actor-Critic method. As shown in Algorithm 2 we extend proximal policy optimization algorithms (PPO) to multi-agent setting by

parameter sharing [24]. At each iteration of the PS-PPO, each agent collects samples according to $\pi_{\theta_k}(o_j)$. After that, we aggregate the samples of all the agents as a batch of data to calculate the advantage function and update network parameters.

---

**Algorithm 2** PS-PPO

---

1: Initialize policy network $\pi_\theta$ and value function $V_\phi$
2: **for** $k = 1$ to $M$ **do**
3:    Agents j=0,1,2,3,...M execute policy $\pi_{\theta_k}(o_j)$,collect trajectories $D_k = \tau_i{}^j$
4:    Trajectories data normalization
5:    Calculate the accumulated discount reward $\widehat{R}_t$
6:    Estimate advantage function $\hat{A}_t$ based on the current value function
7:    Update the policy

$$\theta_{k+1} = \arg\max_\theta \frac{1}{|D_k|T} \sum_{\tau \in D_k} \sum_{t=0}^{T} \min$$

$$\left( \frac{\pi_\theta(a_t|o_t)}{\pi_{\theta_k}(a_t|o_t)} A^{\pi_{\theta_k}}(o_t, a_t), g(\epsilon, A^{\pi_{\theta_k}}(o_t, a_t)) \right)$$

8:    Update the value function

$$\phi_{k+1} = \arg\min_\phi \frac{1}{|D_k|T} \sum_{\tau \in D_k} \sum_{t=0}^{T} (V_\phi(o_t) - \hat{R}_t)^2$$

9: **end for**

---

## III. EXPERIMENTS

We evaluate the performance of the PS-AIRL algorithm in the following two scenarios. The first scenario is conducted in Rendezvous Problem [25], where the goal is to minimize the distances between all the agents. The expert agents were trained by reinforcement learning algorithms using the true reward function. The second scenario is to align each agent's movement direction. The Vicsek model was used to construct the expert policy. We seek to recover both the policy function and the reward function from the expert demonstration. We evaluated PS-AIRL under the setting of different agent numbers as well as different quantities of expert demonstration. Moreover, We compare our methods with behavioral cloning (BC), which learns a maximum likelihood estimate for $a_i$ given each state $s$ and does not require actions from other agents [26].The simulation environment is modified from [25].

### A. Rendezvous Problem

We train the expert agents based on PS-PPO. In this process, the resultant policy $\pi_E$ is constructed to generate the expert demonstration $\tau_E$. Then we use the PS-AIRL to imitate the expert policy and reconstruct the reward function.Imitation performance of different algorithms are evaluated via the true accumulated reward obtained in an episode. Fig. 2 displays the performance of PS-AIRL, Expert, BC and Random policy. PS-AIRL performs consistently better than BC in all the settings. It can be seen from the results that PS-AIRL is not sensitive to the increase of the number of agents. Table I shows the

imitation evaluation results under different number of agents and different number of expert trajectories. It further confirms that PS-AIRL outperforms the other algorithms can achieve a performance that are close to the experts.



Fig. 2. Episode rewards of the imitation policy in the Rendezvous problem. (a) Episode reward with different agent numbers. (b)Episode rewards with different expert demonstrations

We interpret the ability of reward learning of PS-AIRL by visualizing the reward function. We sample an action set from the action space and calculate the reward for each action in a particular scenario. Fig. 2 visualizes the reward function in a time step. Fig. 2(a) is the visualization of the reward function for the focal agent in Fig. 2(b). The horizontal and vertical coordinates are the linear velocity and angular velocity of the agent. The blue color represents the smaller reward value than red color. In the Fig. 2(b), the target agent is driving away from other agents. Fig. 2(a) shows that in the current state, the linear velocity has a small effect on the focal agents' reward, and the focal agent tends to have a higher angular velocity. This observation indicates that the agent will get higher reward value when it turns clockwise or counterclockwise to other agents. This behavior is exactly what is needed to complete the rendezvous goal. The experimental results show that the swarm inverse reinforcement learning algorithm based on PS-AIRL can identify a reasonable reward function.



Fig. 3. Visualization of the learned reward function for the rendezvous agents. (a)The reward under different action for the focal agent.(b)The focal agent is denoted by yellow star.

### B. Vicsek model

We test the PS-AIRL framework on the demonstration data generated by the Vicsek model [2]. At each time instance,

TABLE I
POLICY IMITATING PERFORMANCE IN RENDEZVOUS TASKS

| #Agents | Algorithm | Expert Trajectories | | | |
|---|---|---|---|---|---|
| | | 50 | 80 | 100 | 150 |
| 5 | Expert | -38.86±1.85 | | | |
| | PS-AIRL | -56.78±4.48 | -53.86±3.56 | -52.34±3.34 | -52.28±2.88 |
| | Behavioral cloning | -78.83±16.67 | -75.76±15.67 | -68.68±17.86 | -67.86±14.56 |
| 10 | Expert | -39.74±2.36 | | | |
| | PS-AIRL | -67.86±6.84 | -58.78±5.83 | -53.95±4.78 | -52.35±4.68 |
| | Behavioral cloning | -79.56±23.56 | -78.57±32.58. | -77.07±28.73 | -76.08±27.65 |
| 15 | Expert | -45.86±2.64 | | | |
| | PS-AIRL | -68.24±8.58 | -59.34±6.57 | -57.43±5.87 | -56.38±4.65 |
| | Behavioral cloning | -70.34±24.67 | -77.57±28.78 | -76.87±24.58 | -85.78±23.57 |
| 20 | Expert | -52.76±2.56 | | | |
| | PS-AIRL | -69.32±9.64 | -65.48±8.58 | -63.24±7.69 | -62.76±5.78 |
| | Behavioral cloning | -92.65±18.78 | -89.75±16.69 | -90.56±15.78 | -88.65±13.87 |

the agents' orientations get synchronously updated to the average orientation of their neighbors (including themselves) with additive random perturbations. Our goal is to learn a model for this expert behavior from recorded agent trajectories using the proposed framework. We using the order parameter to evaluate the imitating policy behavior [2].

The overall task performance is evaluated based on their cumulative order parameters. Each episode is fixed 300 time steps. Fig. 4 shows the performance of the expert strategy, PS-AIRL, BC and random strategy. Naturally, the performance of BC increases with more expert demonstrations. It can be seen from the results that PS-AIRL is not sensitive to the increase in the number of agents, and can achieve a performance that are close to the experts. In fact, the parameter sharing mechanism introduced on the basis of the homogeneity assumption makes PS-AIRL is less prone to changes in the number of agents. Table II shows the imitation evaluation results under different number of agents and different number of expert trajectories.



Fig. 4. Episode Order parameter of the imitation policy for the Vicsek model. (a) Episode Order parameter with different agent numbers. (b)Episode Order parameter with different expert demonstrations

The Fig. 5 visualizes the reward function of the Vicsek model. The horizontal and vertical coordinates of the Fig. 5(a) are the linear velocity and angular velocity of the agent. The average motion direction of the agents is shown by the red arrow in Fig. 5(a). There is a certain gap between the direction of the focal agent and the average direction, the focal agent



Fig. 5. Visualization of the learned reward function for the Vicsek model. (a)The reward under different actions for the focal agent.(b) The orientation alignment of ten moving agents. The focal agent is denoted by yellow star.

needs to rotate anticlockwise to meet the motion rules of the Vicsek model. The Fig. 5(a) displays that when the agent rotates anticlockwise, it will obtain a higher reward. Although there is no "true" reward model for the Vicsek system, one can see from the system equations that the agents tend to align over time. The experimental results show that PS-AIRL can reasonably identify the reward function of the Vicsek model.

## IV. DISCUSSION AND FUTURE WORK

We propose a swarm inverse reinforcement learning method called PS-AIRL specifically for collective biological systems. Biological systems are often composed of many agents having large action and observation spaces, challenging for physical models to discover and interpret inner cognition mechanisms. Though the traditional models predict the behaviors, the latent cognitive process for animal collective behaviors is not yet fully understood. PS-AIRL can imitate the expert policy and reconstruct the reward function based on demonstration data. We conduct a theoretical analysis of the parameter-sharing paradigm showing that it can be used to extend the single inverse reinforcement learning to a collective setting in the biological system. Experimental results show that PS-AIRL can achieve excellent performance close to the expert system and reconstruct an explainable reward function for the agents.

Based on the current work on PS-AIRL, we plan to study reward function and strategy function of different specifies

TABLE II
SMALL CAPS: POLICY IMITATING PERFORMANCE FOR VICSEK MODEL

| #Agents | Algorithm | Expert Trajectories | | | |
|---|---|---|---|---|---|
| | | 50 | 80 | 100 | 150 |
| 5 | Expert | 295.32±1.85 | | | |
| | PS-AIRL | 259.06±5.61 | 263.25±5.31 | 266.29±2.46 | 267.30±2.47 |
| | Behavioral cloning | 164.11±31.23 | 201.70±21.72 | 204.51±16.37 | 206.51±13.37 |
| 10 | Expert | 293.53±2.39 | | | |
| | PS-AIRL | 249.96±5.59 | 256.07±5.58 | 263.01±2.67 | 264.03±2.59 |
| | Behavioral cloning | 125.37±10.14 | 207.81±17.72 | 236.59±11.94 | 237.31±10.85 |
| 15 | Expert | 292.41±2.73 | | | |
| | PS-AIRL | 253.46±7.21 | 253.03±4.59 | 257.12±6.20 | 258.13±6.15 |
| | Behavioral cloning | 173.60±10.42 | 233.68±7.62 | 244.55±5.98 | 245.63±5.96 |
| 20 | Expert | 288.42±2.85 | | | |
| | PS-AIRL | 249.54±2.49 | 254.19±5.72 | 260.22±3.87 | 262.77±5.54 |
| | Behavioral cloning | 180.53±15.53 | 235.63±12.62 | 220.55±12.58 | 243.62±13.62 |

on more swarming scenarios. More thorough analysis and interpretation of the inferred reward function should be done to study the relationship between individual interactions and group behavior.

## ACKNOWLEDGMENT

## REFERENCES

[1] I. D. Couzin, J. Krause, R. James, G. D. Ruxton, and N. R. Franks, "Collective memory and spatial sorting in animal groups," Journal of theoretical biology, vol. 218, no. 1, pp. 1-11, 2002.

[2] T. Vicsek, A. Czirók, E. Ben-Jacob, I. Cohen, and O. Shochet, "Novel type of phase transition in a system of self-driven particles," Physical review letters, vol. 75, no. 6, p. 1226, 1995.

[3] F. Cucker and S. Smale, "Emergent behavior in flocks," IEEE Transactions on automatic control, vol. 52, no. 5, pp. 852-862, 2007.

[4] D. J. Sumpter, R. P. Mann, and A. Perna, "The modelling cycle for collective animal behaviour," Interface focus, vol. 2, no. 6, pp. 764-773, 2012.

[5] K. Ried, T. Müller, and H. J. Briegel, "Modelling collective motion based on the principle of agency: General framework and the case of marching locusts," PLoS one, vol. 14, no. 2, p. e0212044, 2019.

[6] P. Torrens, X. Li, and W. A. Griffin, "Building agent-based walking models by machine-learning on diverse databases of space-time trajectory samples," Transactions in GIS, vol. 15, pp. 67-94, 2011. .

[7] F. J. H. Heras, F. Romero-Ferrero, R. C. Hinz, and G. G. de Polavieja, "Deep attention networks reveal the rules of collective motion in zebrafish," PLoS Comput Biol, vol. 15, no. 9, p. e1007354, Sep 2019.

[8] S. Zhou, M. J. Phielipp, J. A. Sefair, S. I. Walker, and H. B. Amor, "Clone swarms: Learning to predict and control multi-robot systems by imitation," in 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2019: IEEE, pp. 4092-4099.

[9] T. Osa, J. Pajarinen, G. Neumann, J. A. Bagnell, P. Abbeel, and J. Peters, "An algorithmic perspective on imitation learning," Foundations and Trends in Robotics, vol. 7, no. 1-2, pp. 1-179, 2018.

[10] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in Proceedings of the twenty-first international conference on Machine learning, 2004, p. 1.

[11] R. Pinsler, M. Maag, O. Arenz, and G. Neumann, "Inverse reinforcement learning of bird flocking behavior," in ICRA Swarms Workshop, 2018.

[12] S. Yamaguchi et al., "Identification of animal behavioral strategies by inverse reinforcement learning," PLoS computational biology, vol. 14, no. 5, p. e1006122, 2018.

[13] G. Papoudakis, F. Christianos, A. Rahman, and S. V. Albrecht, "Dealing with non-stationarity in multi-agent deep reinforcement learning," arXiv preprint arXiv:1906.04737, 2019.

[14] J. Song, H. Ren, D. Sadigh, and S. Ermon, "Multi-agent generative adversarial imitation learning," in Advances in neural information processing systems, 2018, pp. 7461-7472.

[15] L. Yu, J. Song, and S. Ermon, "Multi-agent adversarial inverse reinforcement learning," in International Conference on Machine Learning, 2019: PMLR, pp. 7194-7201.

[16] A. Šošić, W. R. KhudaBukhsh, A. M. Zoubir, and H. Koeppl, "Inverse Reinforcement Learning in Swarm Systems," in Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems, 2017, pp. 1413-1421.

[17] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in Machine learning proceedings 1994: Elsevier, 1994, pp. 157-163.

[18] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in Advances in neural information processing systems, 2000, pp. 1057-1063.

[19] J. K. Gupta, M. Egorov, and M. Kochenderfer, "Cooperative multi-agent control using deep reinforcement learning," in International Conference on Autonomous Agents and Multiagent Systems, 2017: Springer, pp. 66-83.

[20] F. Christianos, G. Papoudakis, A. Rahman, and S. V. Albrecht, "Scaling Multi-Agent Reinforcement Learning with Selective Parameter Sharing," in International Conference on Machine Learning, 2021: PMLR, pp. 1989-1998.

[21] F. Christianos, L. Schäfer, and S. V. Albrecht, "Shared Experience Actor-Critic for Multi-Agent Reinforcement Learning," in Thirty-fourth Conference on Neural Information Processing Systems, 2020: Curran Associates Inc, pp. 10707-10717.

[22] I. J. Goodfellow et al., "Generative Adversarial Networks," arXiv e-prints, p. arXiv:1406.2661. [Online]. Available: https://ui.adsabs.harvard.edu/abs/2014arXiv1406.2661G

[23] J. Fu, K. Luo, and S. Levine, "Learning Robust Rewards with Adversarial Inverse Reinforcement Learning," in International Conference on Learning Representations, 2018.

[24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," p. arXiv:1707.06347. [Online]. Available: https://ui.adsabs.harvard.edu/abs/2017arXiv170706347S

[25] M. Hüttenrauch, S. Adrian, and G. Neumann, "Deep reinforcement learning for swarm systems," Journal of Machine Learning Research, vol. 20, no. 54, pp. 1-31, 2019.

[26] D. A. Pomerleau, "Efficient training of artificial neural networks for autonomous navigation," Neural computation, vol. 3, no. 1, pp. 88-97, 1991.