# Multi-player $H_{\infty}$ Differential Game using On-Policy and Off-Policy Reinforcement Learning

Peiliang An<sup>1</sup>, Mushuang Liu<sup>1</sup>, Yan Wan<sup>1</sup>, and Frank L. Lewis<sup>2</sup>

Abstract—This paper studies a multi-player  $H_\infty$  differential game for systems of general linear dynamics. In this game, multiple players design their control inputs to minimize their cost functions in the presence of worst-case disturbances. We first derive the optimal control and disturbance policies using the solutions to Hamilton-Jacobi-Isaacs (HJI) equations. We then prove that the derived optimal policies stabilize the system and constitute a Nash equilibrium solution. Two integral reinforcement learning (IRL) -based algorithms, including the policy iteration IRL and off-policy IRL, are developed to solve the differential game online. We show that the off-policy IRL can solve the multi-player  $H_\infty$  differential game online without using any system dynamics information. Simulation studies are conducted to validate the theoretical analysis and demonstrate the effectiveness of the developed learning algorithms.

#### I. Introduction

Differential games [1]-[4] have attracted increasing attentions in the control community due to their wide applications in multi-robot systems [5], [6]. Differential games provide a formal mathematical framework to study the coordination, conflict and control of dynamical systems that involve multiple decision-makers (or players) [1]-[4], [7]. Two types of differential games, including the two-player zero-sum games and multi-player nonzero-sum games, have been studied [1], [4]. The two-player zero-sum games can be used to solve the pursuit-evasion type of problems, i.e., there is a single performance index that one player tries to minimize while the other tries to maximize [2], [8]. The two-player zerosum games have also been used to solve the  $H_{\infty}$  control of systems subject to additive external disturbances [1], [8]. The other type of differential games, i.e., the multi-player nonzero-sum games, have been developed to solve the leaderfollower optimal tracking type of problems, where there can generally exist more than two players and each player tries to minimize its individual performance index [3]. In this paper, we study a new type of differential game, called the multiplayer  $H_{\infty}$  differential game, which takes features of the two differential games aforementioned. In the multi-player  $H_{\infty}$ game, each player seeks to minimize its performance index in the presence of a worst-case disturbance. This game provides a theoretical framework for optimal controller design of multi-player systems subject to external disturbances. Per

\*We thank the ONR Grant N00014-18-1-2221, NSF grants 1714519 and 1839804, and ARO Grant W911NF-20-1-0132 for the support of this work.

the knowledge of the authors, there are very limited studies till now that study the multi-player  $H_{\infty}$  differential game [9], [10]. Properties of such systems, e.g., stability and Nash equilibrium have not been thoroughly analyzed.

Finding Nash equilibrium solutions to differential games is not an easy task [3]. In particular, solving zero-sum differential games relies on solving Hamilton-Jacobi-Isaacs (HJI) equations, and solving nonzero-sum differential games relies on solving Hamilton-Jacobi-Bellman (HJB) equations. It has been shown that solving these equations directly in an analytical way is extremely difficult [11]. In addition, solving these equations also requires the information of system dynamics, which is not always available in real applications.

Reinforcement learning (RL) has emerged as an efficient numerical tool for solving optimal control problems online. The use of RL in control theory is documented in [12] for continuous-time linear systems, [13], [14] for discrete-time linear systems, [15], [16] for continuous-time nonlinear systems, and [17] for discrete-time nonlinear systems. Of our interests, RL-based algorithms have also been developed for differential games. Interested readers please refer to [8], [18]–[20] for two-player zero-sum games, and [21], [22] for multi-player nonzero-sum games. In particular, an off-policy integral RL (IRL) was developed in [22] to solve the multiplayer nonzero-sum games without requiring any information of the system dynamics. In this paper, we study both onpolicy and off-policy IRL solutions to the new multi-player  $H_{\infty}$  differential game.

The contributions of this paper are three-fold. First, we formulate the multi-player  $H_{\infty}$  differential game subject to the worst-case external disturbance, and show that the solution to the game stabilizes the system and constitutes a Nash equilibrium. Second, we develop a policy iteration-based learning algorithm to solve the game online, using partial system dynamics information. Third, we further develop an off-policy IRL algorithm that requires no information of the system dynamics.

The remainder of the paper is structured as follows. Section II formulates the multi-player  $H_{\infty}$  differential game and provides preliminaries to facilitate the analysis. In Section III, properties of the multi-player  $H_{\infty}$  game are studied, and two IRL-based algorithms are developed to find the optimal solutions online. Section IV presents simulation studies and Section V concludes the paper.

## II. PROBLEM FORMULATION AND PRELIMINARIES

In this section, we formulate the multi-player  $H_{\infty}$  differential game for a system of general linear dynamics. We then

<sup>&</sup>lt;sup>1</sup>Peiliang An, Mushuang Liu and Yan Wan are with the Department of Electrical Engineering, University of Texas at Arlington, Arlington, TX 76019, USA. Email: peiliang.an@mavs.uta.edu, mushuang.liu@mavs.uta.edu, and yan.wan@uta.edu (contact author).

<sup>&</sup>lt;sup>2</sup> Frank L. Lewis is with UTA Research Institute, University of Texas at Arlington, Fort Worth, Texas, USA. Email: lewis@uta.edu.

provide preliminaries to facilitate the analysis in Section III.

#### A. Problem Formulation

Consider a general N-player linear time-invariant dynamical system given by

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{u}_{j} + \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{d}_{j}, \tag{1}$$

where  $\mathbf{x} = \mathbf{x}(t) \in \mathbf{R}^n$  is the state vector,  $\mathbf{u}_j = \mathbf{u}_j(t) \in \mathbf{R}^m$  is the control input for player j, and the  $\mathbf{d}_j = \mathbf{d}_j(t) \in \mathbf{R}^q$  is the adversarial disturbance input for player j.  $\mathbf{A}$ ,  $\mathbf{B}_j$ , and  $\mathbf{C}_j$  are the drift, control input, and disturbance input dynamics, respectively. It is assumed that the system (1) is stabilizable. Many engineering systems are governed by dynamics (1), for example, the aircraft launching, where  $\mathbf{x}$  is the aircraft speed,  $\mathbf{u}_j$  and  $\mathbf{d}_j$  are the control thrust force and the disturbance force of the controller j, respectively.

Define the cost function to be optimized for player i ( $i = 1, 2, \dots, N$ ) as

$$J_{i}(\mathbf{x}(0), \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i})$$

$$= \int_{0}^{\infty} r_{i}(\mathbf{x}, \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i}) dt$$

$$= \int_{0}^{\infty} \left( \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2} \right) dt,$$
(2)

where  $\mathbf{u}_{-i}$  and  $\mathbf{d}_{-i}$  are the sets of control and disturbance policies for all players other than player i.  $\mathbf{Q}_i$  and  $\mathbf{R}_{ij}$  ( $i \neq j$ ) are positive semi-definite matrices, and  $\mathbf{R}_{ii}$  are positive definite matrices.

The value function of player i is defined as

$$V_{i}(\mathbf{x}(t))$$

$$= \int_{t}^{\infty} r_{i}(\mathbf{x}, \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i}) d\tau$$

$$= \int_{t}^{\infty} \left( \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2} \right) d\tau.$$
(3)

Define the multi-player  $H_{\infty}$  differential game as

$$V_i^*(\mathbf{x}(0)) = \min_{\mathbf{u}_i} \max_{\mathbf{d}_i} J_i(\mathbf{x}(0), \mathbf{u}_i, \mathbf{u}_{-i}, \mathbf{d}_i, \mathbf{d}_{-i}), \quad (4)$$

where  $V_i^*(\mathbf{x}(0))$  is the optimal value for player i. In the multi-player  $H_{\infty}$  game, each player tries to minimize its cost function by choosing a control policy  $\mathbf{u}_i$ , while the disturbance  $\mathbf{d}_i$  seeks to maximize this cost. Each player has access to the full state of the system.

The problem is to find the optimal control and disturbance policies  $\mathbf{u}_{i}^{*}$  and  $\mathbf{d}_{i}^{*}$  such that

$$\mathbf{u}_{i}^{*} = \underset{\mathbf{u}_{i}}{\operatorname{argmin}} J_{i}(\mathbf{x}(0), \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i}),$$
  
$$\mathbf{d}_{i}^{*} = \underset{\mathbf{d}_{i}}{\operatorname{argmax}} J_{i}(\mathbf{x}(0), \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i}).$$

#### B. Preliminaries

**Definition 1.** [11] The system (1) is said to have  $L_2$ -gain less than or equal to  $\gamma$  if the following disturbance attenuation condition is satisfied for all  $\mathbf{d}_i \in L_2[0,\infty)$  with  $\mathbf{x}(0) = \mathbf{0}$ :

$$\frac{\int_t^\infty \|\mathbf{z}(\tau)\|^2 d\tau}{\int_t^\infty \left(\sum_{j=1}^N \|\mathbf{d}_j\|^2\right) d\tau} \le \gamma^2,$$

where  $\|\mathbf{z}(t)\|^2 = \mathbf{x}^{\mathsf{T}} \mathbf{Q}_i \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_j^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_j$ ,  $\mathbf{d}_j(t)$  is the disturbance input, and  $\gamma$  is the amount of attenuation.

It is assumed that  $\gamma$  in (2) satisfies  $\gamma \geq \gamma^*$ , where  $\gamma^*$  is the smallest  $\gamma$ , also know as  $H_{\infty}$  gain for system (1) [1], which satisfies the disturbance attenuation condition.

**Definition 2.** [1] Policies  $\{\mathbf{u}_1^*, \mathbf{d}_1^*, \mathbf{u}_2^*, \mathbf{d}_2^*, \cdots, \mathbf{u}_N^*, \mathbf{d}_N^*\}$  are said to constitute a Nash equilibrium solution to the N-player  $H_{\infty}$  game if the following inequality holds:

$$J_{i}(\mathbf{x}(0), \mathbf{u}_{i}^{*}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}, \mathbf{d}_{-i}^{*})$$

$$\leq J_{i}^{*}(\mathbf{x}(0), \mathbf{u}_{i}^{*}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}^{*}, \mathbf{d}_{-i}^{*})$$

$$\leq J_{i}(\mathbf{x}(0), \mathbf{u}_{i}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}^{*}, \mathbf{d}_{-i}^{*}), \quad \forall \mathbf{u}_{i}, \forall \mathbf{d}_{i}, \forall i.$$
(5)

# III. MULTI-PLAYER $H_{\infty}$ DIFFERENTIAL GAME

This section derives the optimal solution to the N-player  $H_{\infty}$  differential game. Section III-A studies the stability and Nash equilibrium of the game. Two IRL-based algorithms are then developed in Section III-B to solve the differential game online.

# A. Stability and Nash Equilibrium

Differentiating the value function  $V_i(\mathbf{x}(t))$  defined in (3), one can obtain the Bellman equation as follows,

$$\mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2}$$

$$+ \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} + \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{u}_{j} + \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{d}_{j} \right) = 0,$$
(6)

where  $\nabla V_i = \partial V_i/\partial \mathbf{x}$ . The boundary condition for this partial differential equation is  $V_i(\mathbf{0}) = 0$ . A solution to (6) is the value function  $V_i(\mathbf{x})$  for the feedback control policy  $\mathbf{u}_i = \mathbf{u}_i(V_i(\mathbf{x}))$  and disturbance policy  $\mathbf{d}_i = \mathbf{d}_i(V_i(\mathbf{x}))$ .

Define the Hamiltonian function associated with the value function (3) as

$$H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i})$$

$$= r_{i}(\mathbf{x}, \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i})$$

$$+ \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} + \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{u}_{j} + \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{d}_{j} \right)$$

$$= \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2}$$

$$+ \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} + \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{u}_{j} + \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{d}_{j} \right).$$

$$(7)$$

At the equilibrium point, applying the stationary conditions

$$\frac{\partial H_i}{\partial \mathbf{u}_i} = 0$$
 and  $\frac{\partial H_i}{\partial \mathbf{d}_i} = 0$ 

yields the optimal control and disturbance policies as functions of  $V_i(\mathbf{x})$ :

$$u_i^* = \mathbf{u}_i^*(V_i(\mathbf{x})) = -\frac{1}{2}\mathbf{R}_{ii}^{-1}\mathbf{B}_i^{\mathsf{T}}\nabla V_i, \tag{8}$$

$$\mathbf{d}_{i}^{*} = \mathbf{d}_{i}^{*}(V_{i}(\mathbf{x})) = \frac{1}{2\gamma^{2}} \mathbf{C}_{i}^{\mathsf{T}} \nabla V_{i}. \tag{9}$$

Therefore, the value function  $V_i(\mathbf{x})$  in (3) is only a function of the state  $\mathbf{x}(t)$ . Moreover, the Hamiltonian function  $H_i$  attains a saddle point at the stationary point since  $\partial^2 H_i/\partial \mathbf{u}_i^2 = 2\mathbf{R}_{ii} > 0$  and  $\partial^2 H_i/\partial \mathbf{d}_i^2 = -2\gamma^2 < 0$ .

Substituting (8) and (9) into the Bellman Equation (6), the following Hamilton-Jacobi-Isaacs (HJI) equation is obtained:

$$\mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \frac{1}{4} \sum_{j=1}^{N} \nabla V_{j}^{\mathsf{T}} \mathbf{B}_{j} \mathbf{R}_{jj}^{-1} \mathbf{R}_{ij} \mathbf{R}_{jj}^{-1} \mathbf{B}_{j}^{\mathsf{T}} \nabla V_{j}$$

$$- \frac{1}{4\gamma^{2}} \sum_{j=1}^{N} \nabla V_{j}^{\mathsf{T}} \mathbf{C}_{j} \mathbf{C}_{j}^{\mathsf{T}} \nabla V_{j} + \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} \right)$$

$$- \frac{1}{2} \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{R}_{jj}^{-1} \mathbf{B}_{j}^{\mathsf{T}} \nabla V_{j} + \frac{1}{2\gamma^{2}} \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{C}_{j}^{\mathsf{T}} \nabla V_{j} \right) = 0.$$
(10)

Since the attenuation condition in Definition 1 is satisfied, the HJI equation (10) has a positive semi-definite solution  $V_i^*(\mathbf{x}(t))$  [1].

Note that for the optimal policies  $\mathbf{u}_i^*$ ,  $\mathbf{d}_i^*$  and the corresponding  $V_i^*$ , the HJI equation satisfies

$$H_i(\mathbf{x}, \nabla V_i^*, \mathbf{u}_i^*, \mathbf{u}_{-i}^*, \mathbf{d}_i^*, \mathbf{d}_{-i}^*) = 0.$$
 (11)

**Theorem 1.** Assume the control and disturbance policies are optimal for all players other than player i. Then for any admissible policies  $\mathbf{u}_i(\mathbf{x})$  and  $\mathbf{d}_i(\mathbf{x})$ , and any positive semi-definite value function  $V_i(\mathbf{x})$ , one has the following equation:

$$H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}, \mathbf{d}_{-i}^{*})$$

$$= H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}^{*}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}^{*}, \mathbf{d}_{-i}^{*})$$

$$+ (\mathbf{u}_{i} - \mathbf{u}_{i}^{*})^{\mathsf{T}} \mathbf{R}_{ii} (\mathbf{u}_{i} - \mathbf{u}_{i}^{*}) - \gamma^{2} (\mathbf{d}_{i} - \mathbf{d}_{i}^{*})^{\mathsf{T}} (\mathbf{d}_{i} - \mathbf{d}_{i}^{*}).$$
(12)

*Proof*: Taking  $\mathbf{u}_{-i} = \mathbf{u}_{-i}^*$  and  $\mathbf{d}_{-i} = \mathbf{d}_{-i}^*$ , the Hamiltonian function in (7) can be written as

$$H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}, \mathbf{d}_{-i}^{*})$$

$$= \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j \neq i} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j}^{*} + \mathbf{u}_{i}^{\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i}$$

$$- \gamma^{2} \sum_{j \neq i} \|\mathbf{d}_{j}^{*}\|^{2} - \gamma^{2} \mathbf{d}_{i}^{\mathsf{T}} \mathbf{d}_{i}$$

$$+ \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} + \sum_{j \neq i} \mathbf{B}_{j} \mathbf{u}_{j}^{*} + \mathbf{B}_{i} \mathbf{u}_{i} + \sum_{j \neq i} \mathbf{C}_{j} \mathbf{d}_{j}^{*} + \mathbf{C}_{i} \mathbf{d}_{i} \right)$$

$$= \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j} \mathbf{u}_{j}^{*\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j}^{*} - \gamma^{2} \sum_{j} \|\mathbf{d}_{j}^{*}\|^{2}$$

$$+ \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} + \sum_{j} \mathbf{B}_{j} \mathbf{u}_{j}^{*} + \sum_{j} \mathbf{C}_{j} \mathbf{d}_{j}^{*} \right)$$

$$+ \mathbf{u}_{i}^{\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i} - \mathbf{u}_{i}^{*\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i}^{*} - \gamma^{2} \mathbf{d}_{i}^{\mathsf{T}} \mathbf{d}_{i} + \gamma^{2} \mathbf{d}_{i}^{*\mathsf{T}} \mathbf{d}_{i}^{*}$$

$$+ \nabla V_{i}^{\mathsf{T}} (\mathbf{B}_{i} \mathbf{u}_{i} - \mathbf{B}_{i} \mathbf{u}_{i}^{*} + \mathbf{C}_{i} \mathbf{d}_{i} - \mathbf{C} \mathbf{d}_{i}^{*})$$

$$= H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}^{*}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}^{*}, \mathbf{d}_{-i}^{*}) + \mathbf{u}_{i}^{\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i} - \mathbf{u}_{i}^{*\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i}^{*}$$

$$- \gamma^{2} \mathbf{d}_{i}^{\mathsf{T}} \mathbf{d}_{i} + \gamma^{2} \mathbf{d}_{i}^{*\mathsf{T}} \mathbf{d}_{i}^{*} + (\mathbf{u}_{i} - \mathbf{u}_{i}^{*})^{\mathsf{T}} \mathbf{B}_{i}^{\mathsf{T}} \nabla V_{i}$$

$$+ (\mathbf{d}_{i} - \mathbf{d}_{i}^{*})^{\mathsf{T}} \mathbf{C}_{i}^{\mathsf{T}} \nabla V_{i}.$$
(13)

According to (8) and (9), one has

$$\mathbf{B}_{i}^{\mathsf{T}} \nabla V_{i} = -2\mathbf{R}_{ii} \mathbf{u}_{i}^{*}$$
 and  $\mathbf{C}_{i}^{\mathsf{T}} \nabla V_{i} = 2\gamma^{2} \mathbf{d}_{i}^{*}$ .

As such, (13) can be further rewritten as

$$H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}, \mathbf{d}_{-i}^{*})$$

$$= H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}^{*}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}^{*}, \mathbf{d}_{-i}^{*}) + \mathbf{u}_{i}^{\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i}$$

$$- \mathbf{u}_{i}^{\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i}^{*} - \gamma^{2} \mathbf{d}_{i}^{\mathsf{T}} \mathbf{d}_{i} + \gamma^{2} \mathbf{d}_{i}^{\mathsf{T}} \mathbf{d}_{i}^{*}$$

$$- 2(\mathbf{u}_{i} - \mathbf{u}_{i}^{*})^{\mathsf{T}} \mathbf{R}_{ii} \mathbf{u}_{i}^{*} + 2\gamma^{2} (\mathbf{d}_{i} - \mathbf{d}_{i}^{*})^{\mathsf{T}} \mathbf{d}_{i}^{*}$$

$$= H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}^{*}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}^{*}, \mathbf{d}_{-i}^{*})$$

$$+ (\mathbf{u}_{i} - \mathbf{u}_{i}^{*})^{\mathsf{T}} \mathbf{R}_{ii} (\mathbf{u}_{i} - \mathbf{u}_{i}^{*}) - \gamma^{2} (\mathbf{d}_{i} - \mathbf{d}_{i}^{*})^{\mathsf{T}} (\mathbf{d}_{i} - \mathbf{d}_{i}^{*}).$$

This result is next employed to show that the optimal policies given by (8) and (9) in terms of coupled HJI solution  $V_i^*(\mathbf{x})$  constitute a Nash equilibrium solution.

**Theorem 2.** Suppose  $V_i^*(\mathbf{x})$  are smooth continuous positive semi-definite functions that solve the HJI equations (10). The control and disturbance policies  $\mathbf{u}_i^*$  and  $\mathbf{d}_i^*$  are given by (8) and (9). Then the following two statements (a) and (b) hold. (a). The closed-loop system

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{u}_{j}^{*} + \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{d}_{j}^{*}$$

$$= \mathbf{A}\mathbf{x} - \frac{1}{2} \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{R}_{jj}^{-1} \mathbf{B}_{j}^{\mathsf{T}} \nabla V_{j}^{*} + \frac{1}{2\gamma^{2}} \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{C}_{j}^{\mathsf{T}} \nabla V_{j}^{*}$$
(14)

is asymptotically stable.

(b). Policies  $\{u_i^*, d_i^*\}$  constitute a Nash solution.

Proof:

(a). With  $\gamma$  satisfying the attenuation condition, one has

$$= \int_{t}^{\infty} \left( \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2} \right) d\tau \geq 0,$$

where  $V_i(\mathbf{x}) = 0$  if and only if  $\mathbf{x} = \mathbf{0}$ .

Select  $V_i(\mathbf{x})$  as the Lyapunov function candidates. Differentiating  $V_i(\mathbf{x})$  yields

$$\begin{split} \dot{V}_i(\mathbf{x}) &= \left(\nabla V_i\right)^{^{\mathrm{T}}} \left(\mathbf{A}\mathbf{x} + \sum_{j=1}^N \mathbf{B}_j \mathbf{u}_j + \sum_{j=1}^N \mathbf{C}_j \mathbf{d}_j\right) \\ &= -\left(\mathbf{x}^{^{\mathrm{T}}} \mathbf{Q}_i \mathbf{x} + \sum_{j=1}^N \mathbf{u}_j^{^{\mathrm{T}}} \mathbf{R}_{ij} \mathbf{u}_j - \gamma^2 \sum_{j=1}^N \|\mathbf{d}_j\|^2\right) \leq 0, \end{split}$$

where  $\dot{V}_i(\mathbf{x}) = 0$  if and only if  $\mathbf{x} = \mathbf{0}$ . Therefore,  $V_i(\mathbf{x})$  are Lynapunov functions and the system (14) is asymptotically stable.

(b). Since the system (14) is asymptotically stable, one has  $\mathbf{x}(t) \to \mathbf{0}$ , and thus  $V_i(\mathbf{x}(t)) \to 0$ , as time  $t \to \infty$ . The cost function (2) can be rewritten as

$$J_{i}(\mathbf{x}(0), \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i})$$

$$= \int_{0}^{\infty} \left( \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2} \right) dt$$

$$+ \int_{0}^{\infty} \dot{V}_{i} dt - V_{i}(\mathbf{x}(\infty)) + V_{i}(\mathbf{x}(0))$$

$$= \int_{0}^{\infty} \left( \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{\mathsf{T}} \mathbf{R}_{ij} \mathbf{u}_{j} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}\|^{2} \right) dt$$

$$+ \int_{0}^{\infty} \nabla V_{i}^{\mathsf{T}} \left( \mathbf{A} \mathbf{x} + \sum_{j=1}^{N} \mathbf{B}_{j} \mathbf{u}_{j} + \sum_{j=1}^{N} \mathbf{C}_{j} \mathbf{d}_{j} \right) dt$$

$$+ V_{i}(\mathbf{x}(0))$$

$$= \int_{0}^{\infty} H_{i}(\mathbf{x}, \nabla V_{i}, \mathbf{u}_{i}, \mathbf{u}_{-i}, \mathbf{d}_{i}, \mathbf{d}_{-i}) dt + V_{i}(\mathbf{x}(0)).$$

Now let  $V_i(\mathbf{x}) = V_i^*(\mathbf{x})$  satisfy the HJI equation (10), and  $\mathbf{u}_{-i}$ ,  $\mathbf{d}_{-i}$  choose the optimal policies. By Theorem 1 one has

$$J_{i}(\mathbf{x}(0), \mathbf{u}_{i}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}, \mathbf{d}_{-i}^{*})$$

$$= \int_{0}^{\infty} H_{i}(\mathbf{x}, \nabla V_{i}^{*}, \mathbf{u}_{i}, \mathbf{u}_{-i}^{*}, \mathbf{d}_{i}, \mathbf{d}_{-i}^{*}) dt + V_{i}^{*}(\mathbf{x}(0))$$

$$= \int_{0}^{\infty} \left( (\mathbf{u}_{i} - \mathbf{u}_{i}^{*})^{\mathsf{T}} \mathbf{R}_{ii} (\mathbf{u}_{i} - \mathbf{u}_{i}^{*}) - \gamma^{2} (\mathbf{d}_{i} - \mathbf{d}_{i}^{*})^{\mathsf{T}} \right)$$

$$(\mathbf{d}_{i} - \mathbf{d}_{i}^{*}) dt + V_{i}^{*}(\mathbf{x}(0)),$$

which implies that (5) is satisfied and hence the system is in Nash equilibrium.

## B. Approximated Solutions Using IRL

In Section III-A, we develop the optimal policies for the multi-player  $H_{\infty}$  differential game. As one may notice, the key to finding the policies is solving  $V_i^*(x)$  from the HJI Equation (10), which is, however, extremely difficult analytically [3]. As such, we propose two IRL-based algorithms to solve the HJI equation numerically.

1) On-Policy IRL: The value function (3) can be written

$$V_{i}(\mathbf{x}(t)) = \int_{t}^{t+T} r_{i}(\mathbf{x}(\tau), \mathbf{u}_{i}(\tau), \mathbf{u}_{-i}(\tau), \mathbf{d}_{i}(\tau), \mathbf{u}_{-i}(\tau)) d\tau$$
(15)  
+  $V_{i}(\mathbf{x}(t+T)),$ 

where T is the time interval. Assume that there exits a weight vector  $\mathbf{W}_i$  such that the value function can be approximated as

$$V_i(\mathbf{x}) = \mathbf{W}_i^{\mathsf{T}} \phi_i(\mathbf{x}), \tag{16}$$

where  $\phi_i(\mathbf{x})$  is the basis function vector.

With the approximated value function, the optimal control and disturbance policies can then be determined using RL, in particular, the Policy Iteration (PI) algorithm [1, Page 474]. The PI algorithm constitutes two iterative steps: Policy Evaluation step, to evaluate the value function by (15) and (16), and Policy Improvement step, to find the optimal policies based on current value function by (8) and (9). This PI algorithm for the multi-player  $H_{\infty}$  differential game is summarized in Algorithm 1.

**Algorithm 1** Policy iteration algorithm for multi-player  $H_{\infty}$  differential game

- 1: Initialize each player with admissible policies  $\mathbf{u}_i^{(1)}$  and  $\mathbf{d}_i^{(1)}$ ,  $\forall i \in N$ .
- 2: For each iteration k, find the value function  $V_i^{(k)}(t)$  by

$$V_i^{(k)}(\mathbf{x}(t)) = \int_t^{t+T} r_i \left( \mathbf{x}, \mathbf{u}_i^{(k)}, \mathbf{u}_{-i}^{(k)}, \mathbf{d}_i^{(k)}, \mathbf{d}_{-i}^{(k)} \right) d\tau$$
$$+ \mathbf{W}_i^{(k-1)^{\mathsf{T}}} \phi_i(\mathbf{x}(t+T)). \tag{17}$$

3: Update the weight vector  $\mathbf{W}_i^{(k)}$  according to the estimated  $V_i^{(k)}(\mathbf{x}(t))$  using the least-squares method,

$$\mathbf{W}_{i}^{(k)^{\mathsf{T}}} \phi_{i}(\mathbf{x}(t)) = V_{i}^{(k)}(\mathbf{x}(t)). \tag{18}$$

4: Update the policies  $\mathbf{u}_{i}^{(k+1)}$  and  $\mathbf{d}_{i}^{(k+1)}$  for all players as

$$\mathbf{u}_{i}^{(k+1)} = -\frac{1}{2} \mathbf{R}_{ii}^{-1} \mathbf{B}_{i}^{\mathsf{T}} \frac{\partial V_{i}^{(k)}}{\partial \mathbf{x}},$$

$$\mathbf{d}_{i}^{(k+1)} = \frac{1}{2\gamma^{2}} \mathbf{C}_{i}^{\mathsf{T}} \frac{\partial V_{i}^{(k)}}{\partial \mathbf{x}}.$$
(19)

- 5: Repeat procedures 2-4 until convergence.
- 2) Off-policy IRL: The on-policy algorithm requires the knowledge of the system dynamics, i.e., matrices  $\mathbf{B}_i$  and  $\mathbf{C}_i$ , for learning the optimal policies. In addition, the behavior policies  $\mathbf{u}_i$  and  $\mathbf{d}_i$  are required to be adjustable at every policy improvement step.

This subsection develops an off-policy IRL algorithm to learn the optimal policies without any information of the system dynamics. The off-policy IRL learns the optimal policies of the game online while the game is being played based on fixed behavior policies  $\mathbf{u}_i$  and  $\mathbf{d}_i$ , which are used

to generate system data [11]. This result is developed for the case when players have identical dynamics, i.e.,  $B_i = B$ and  $C_i = C$ , for all  $j = 1, 2, \dots, N$ .

We write the system dynamics in the following form:

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \sum_{j=1}^{N} \mathbf{B}\mathbf{u}_{j}^{(k)} + \sum_{j=1}^{N} \mathbf{C}\mathbf{d}_{j}^{(k)} + \sum_{j=1}^{N} \mathbf{B}\left(\mathbf{u}_{j} - \mathbf{u}_{j}^{(k)}\right) + \sum_{j=1}^{N} \mathbf{C}\left(\mathbf{d}_{j} - \mathbf{d}_{j}^{(k)}\right),$$
(20)

where  $\mathbf{u}_{j}^{(k)}$  and  $\mathbf{d}_{j}^{(k)}$  are the policies to be updated for the

Differentiation the value  $V_i^{(k)}(\mathbf{x}(t))$  along with the system dynamics (20) and using (6), (19) yield

$$\dot{V}_{i}^{(k)}(\mathbf{x}(t))$$

$$= \nabla V_{i}^{(k)^{\mathsf{T}}} \left( \mathbf{A} \mathbf{x} + \sum_{j=1}^{N} \mathbf{B} \mathbf{u}_{j}^{(k)} + \sum_{j=1}^{N} \mathbf{C} \mathbf{d}_{j}^{(k)} \right)$$

$$+ \nabla V_{i}^{(k)^{\mathsf{T}}} \left( \sum_{j=1}^{N} \mathbf{B} \left( \mathbf{u}_{j} - \mathbf{u}_{j}^{(k)} \right) + \sum_{j=1}^{N} \mathbf{C} \left( \mathbf{d}_{j} - \mathbf{d}_{j}^{(k)} \right) \right)$$

$$= - \left( \mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} + \sum_{j=1}^{N} \mathbf{u}_{j}^{(k)^{\mathsf{T}}} \mathbf{R}_{ij} \mathbf{u}_{j}^{(k)} - \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}^{(k)}\|^{2} \right)$$

$$- 2 \mathbf{u}_{i}^{(k+1)^{\mathsf{T}}} \mathbf{R}_{ii} \sum_{j=1}^{N} \left( \mathbf{u}_{j} - \mathbf{u}_{j}^{(k)} \right)$$

$$+ 2 \gamma^{2} \mathbf{d}_{i}^{(k+1)^{\mathsf{T}}} \sum_{j=1}^{N} \left( \mathbf{d}_{j} - \mathbf{d}_{j}^{(k)} \right).$$
(21)

Integrating (21) from both sides gives the following offpolicy IRL Bellman equation:

$$V_{i}^{(k)}(\mathbf{x}(t+T)) - V_{i}^{(k)}(\mathbf{x}(t))$$

$$= \int_{t}^{t+T} \left( -\mathbf{x}^{\mathsf{T}} \mathbf{Q}_{i} \mathbf{x} - \sum_{j=1}^{N} \mathbf{u}_{j}^{(k)^{\mathsf{T}}} \mathbf{R}_{ij} \mathbf{u}_{j}^{(k)} + \gamma^{2} \sum_{j=1}^{N} \|\mathbf{d}_{j}^{(k)}\|^{2} \right) d\tau$$

$$+ \int_{t}^{t+T} \left( -2\mathbf{u}_{i}^{(k+1)^{\mathsf{T}}} \mathbf{R}_{ii} \sum_{j=1}^{N} \left( \mathbf{u}_{j} - \mathbf{u}_{j}^{(k)} \right) \right)$$

$$+ 2\gamma^{2} \mathbf{d}_{i}^{(k+1)^{\mathsf{T}}} \sum_{j=1}^{N} \left( \mathbf{d}_{j} - \mathbf{d}_{j}^{(k)} \right) d\tau.$$
(22)

Note that for any fixed admissible control and disturbance policies  $\mathbf{u}_i$  and  $\mathbf{d}_i$ , (22) can be solved for value function  $V_i^{(k)}$  and the updated policies  $\mathbf{u}_i^{(k+1)}$  and  $\mathbf{d}_i^{(k+1)}$  simultaneously. To this end, three neural networks (NNs), i.e., the critic NN, the actor NN, and the disturber NN, are used here for approximating the value function and the updated control and disturbance policies respectively:

$$V_i^{(k)}(\mathbf{x}) = \mathbf{W}_i^{(k)} \phi_i(\mathbf{x}),$$

$$\mathbf{u}_i^{(k+1)}(\mathbf{x}) = \mathbf{W}_{u,i}^{(k+1)} \sigma_i(\mathbf{x}),$$

$$\mathbf{d}_i^{(k+1)}(\mathbf{x}) = \mathbf{W}_{d,i}^{(k+1)} \psi_i(\mathbf{x}),$$
(23)

where  $\phi_i(\mathbf{x})$ ,  $\sigma_i(\mathbf{x})$  and  $\psi_i(\mathbf{x})$  provide suitable basis function vectors, and  $\mathbf{W}_i^{(k)}$ ,  $\mathbf{W}_{u,i}^{(k+1)}$  and  $\mathbf{W}_{d,i}^{(k+1)}$  are weight matrices with proper dimensions.

The implementation of the off-policy IRL algorithm is described in Algorithm 2.

Algorithm 2 Off-policy IRL algorithm for multi-player  $H_{\infty}$  differential game

- 1: Initialize each player with admissible policies  $\mathbf{u}_{i}^{(1)}$  and
- 2: For each iteration k, solve (22) for  $V_i^{(k)}$ ,  $\mathbf{u}_i^{(k+1)}$ , and
- $\mathbf{d}_{i}^{(k+1)} \text{ simultaneously.}$ 3: Update  $\mathbf{W}_{i}^{(k)}$ ,  $\mathbf{W}_{u,i}^{(k+1)}$  and  $\mathbf{W}_{d,i}^{(k+1)}$  according to the derived  $V_{i}^{(k)}$ ,  $\mathbf{u}_{i}^{(k+1)}$ ,  $\mathbf{d}_{i}^{(k+1)}$  by (23) using the least-
- 4: Repeat procedures 2-3 until convergence

## IV. ILLUSTRATIVE EXAMPLES

In this section, the two proposed algorithms are applied to a linear system example to validate the theoretical analysis.

Consider a three-player  $H_{\infty}$  game with a linear system described by the following dynamics:

$$\dot{\mathbf{x}} = \begin{bmatrix} 1 & 0.25 \\ 1 & 0 \end{bmatrix} \mathbf{x} + \sum_{j}^{3} \begin{bmatrix} 1.3 \\ 0 \end{bmatrix} \mathbf{u}_{j} + \sum_{j}^{3} \begin{bmatrix} 1.3 \\ 0 \end{bmatrix} \mathbf{d}_{j}, \quad (24)$$

where  $\mathbf{x} = [x_1, x_2]^{T}$ .

The parameters in the value function (3) are selected as:  $\mathbf{Q}_1 = \mathbf{Q}_2 = \mathbf{Q}_3 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \ \mathbf{R}_{12} = \mathbf{R}_{13} = \mathbf{R}_{21} = \mathbf{R}_{23} = \mathbf{R}_{31} = \mathbf{R}_{32} = 1, \ \mathbf{R}_{11} = 2, \ \mathbf{R}_{22} = 3, \ \mathbf{R}_{33} = 5, \ \text{and} \ \gamma = 5.$ The reinforcement learning interval T is chosen to be 0.1.

The on-policy PI algorithm (Algorithm 1) is implemented first. We select the basis function  $\phi_i = [x_1^2, x_1x_2, x_2^2]^{^{\mathrm{T}}}$  with weight vector  $\mathbf{W}_i = [W_{i1}, W_{i2}, W_{i3}]^{^{\mathrm{T}}}$ , where i = 1, 2, 3. Figure 1(a) and 1(b) show the evolution of the system states and value function weights.

Figure 1(a) shows that the system states converge to 0 when the optimal policies are applied to the system (24). Moreover, Figure 1(b) verifies the convergence of value function weights, from which the optimal policies can be derived.

Then we simulate the off-policy IRL algorithm (Algorithm 2). Here, three NNs are selected as follows: the critic NN  $\phi_i = [x_1^2, x_1 x_2, x_2^2]^{^{\mathrm{T}}}$  with a weight vector  $\mathbf{W}_i = [W_{i1}, W_{i2}, W_{i3}]^{^{\mathrm{T}}}$ ; the actor NN  $\sigma_i = [x_1, x_2]^{^{\mathrm{T}}}$  with a weight vector  $\mathbf{W}_{u,i} = [W_{u,i1}, W_{u,i2}]^{\mathrm{T}}$ ; the disturber  $NN_{u,i2} = V_{u,i1}$  $[x_1, x_2]^{^{\mathrm{T}}}$  with a weight vector  $\mathbf{W}_{d,i} = [W_{d,i1}, W_{d,i2}]^{^{\mathrm{T}}}$ , where

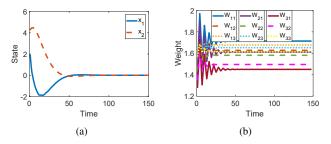


Fig. 1. Multi-player  $H_{\infty}$  differential game using on-policy IRL. (a) The evolution of the system states, and (b) the derived value function weights.

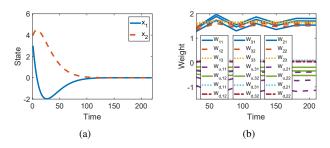


Fig. 2. Multi-player  $H_{\infty}$  differential game using off-policy IRL. (a) The evolution of the system states, and (b) the derived value function weights.

i = 1, 2, 3. The simulation results are shown in Figure 2(a) and Figure 2(b).

Figure 2 shows that the value function weights converge in limited time using the proposed off-policy IRL algorithm, and the converged values are identical to the ones derived from the on-policy algorithm. In addition, the system states converge to **0**, which validate the asymptotic stability of the closed-loop system.In addition, we find that the HJI Equation (10) holds after substituting the derived value function, which verifies the correctness of the derived solutions (18), (19) and (23).

### V. CONCLUSION

This paper studies a new differential game that takes features of two existing games, i.e., two-player zero-sum and multi-player nonzero-sum games, to solve the optimal control problems of multi-player systems subject to external disturbances. We showed that the optimal solutions to this differential game can be found by solving the HJI equation, and the derived optimal solutions can make the system asymptotically stable and in Nash equilibrium. Moreover, to solve the differential games online, we designed two IRL-based algorithms, including the policy iteration and off-policy IRLs. In particular, the designed off-policy IRL can find the Nash solutions without using any information of the system dynamics. In the future, we will generalize the current work to systems with general nonlinear dynamics, and apply the designed algorithms in real-world applications.

# REFERENCES

[1] F. L. Lewis, D. Vrabie, and V. L. Syrmos, *Optimal control*. John Wiley & Sons, 2012.

- [2] T. Başar and P. Bernhard,  $H_{\infty}$  optimal control and related minimax design problems: a dynamic game approach. Springer Science & Business Media, 2008.
- [3] K. G. Vamvoudakis, H. Modares, B. Kiumarsi, and F. L. Lewis, "Game theory-based control system algorithms with real-time reinforcement learning: How to solve multiplayer games online," *IEEE Control Systems*, vol. 37, no. 1, pp. 33–52, 2017.
- [4] M. Liu, Y. Wan, F. Lewis, and V. G. Lopez, "Adaptive optimal control for stochastic multi-player differential games using on-policy and offpolicy reinforcement learning," accepted by IEEE Transactions on Neural Network and Learning Systems, 2020.
- [5] A. Perelman, T. Shima, and I. Rusnak, "Cooperative differential games strategies for active aircraft protection from a homing missile," *Journal* of Guidance, Control, and Dynamics, vol. 34, no. 3, pp. 761–773, 2011
- [6] T. Mylvaganam, M. Sassano, and A. Astolfi, "A differential game approach to multi-agent collision avoidance," *IEEE Transactions on Automatic Control*, vol. 62, no. 8, pp. 4229–4235, 2017.
- [7] A. W. Starr and Y.-C. Ho, "Nonzero-sum differential games," *Journal of optimization theory and applications*, vol. 3, no. 3, pp. 184–206, 1969
- [8] M. Liu, Y. Wan, F. L. Lewis, and V. G. Lopez, "Stochastic twoplayer zero-sum learning differential games," in *Proceedings of IEEE* 15th International Conference on Control and Automation (ICCA), Edinburgh, Scotland, 2019.
- [9] R. Song, Q. Wei, and B. Song, "Neural-network-based synchronous iteration learning method for multi-player zero-sum games," *Neuro-computing*, vol. 242, pp. 73–82, 2017.
- [10] H. Jiang, H. Zhang, J. Han, and K. Zhang, "Iterative adaptive dynamic programming methods with neural network implementation for multiplayer zero-sum games," *Neurocomputing*, vol. 307, pp. 54–60, 2018.
- [11] H. Modares, F. L. Lewis, and Z.-P. Jiang, "h<sub>∞</sub> tracking control of completely unknown continuous-time systems via off-policy reinforcement learning." *IEEE Transactions on Neural Networks and Learning* Systems, vol. 26, no. 10, pp. 2550–2562, 2015.
- [12] D. Vrabie, O. Pastravanu, M. Abu-Khalaf, and F. L. Lewis, "Adaptive optimal control for continuous-time linear systems based on policy iteration," *Automatica*, vol. 45, no. 2, pp. 477–484, 2009.
- [13] F. L. Lewis and D. Vrabie, "Reinforcement learning and adaptive dynamic programming for feedback control," *IEEE Circuits and Systems Magazine*, vol. 9, no. 3, 2009.
- [14] B. Kiumarsi, F. L. Lewis, H. Modares, A. Karimpour, and M.-B. Naghibi-Sistani, "Reinforcement q-learning for optimal tracking control of linear discrete-time systems with unknown dynamics," *Automatica*, vol. 50, no. 4, pp. 1167–1175, 2014.
- [15] K. G. Vamvoudakis and F. L. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878–888, 2010.
  [16] D. Vrabie and F. Lewis, "Neural network approach to continuous-
- [16] D. Vrabie and F. Lewis, "Neural network approach to continuoustime direct adaptive optimal control for partially unknown nonlinear systems," *Neural Networks*, vol. 22, no. 3, pp. 237–246, 2009.
- [17] B. Kiumarsi and F. L. Lewis, "Actor-critic-based optimal tracking for partially unknown nonlinear discrete-time systems," *IEEE Transac*tions on Neural Networks and Learning Systems, vol. 26, no. 1, pp. 140–151, 2015.
- [18] D. Vrabie and F. Lewis, "Adaptive dynamic programming for online solution of a zero-sum differential game," *Journal of Control Theory* and Applications, vol. 9, no. 3, pp. 353–360, 2011.
- [19] H.-N. Wu and B. Luo, "Simultaneous policy update algorithms for learning the solution of linear continuous-time  $h_{\infty}$  state feedback control," *Information Sciences*, vol. 222, pp. 472–485, 2013.
- [20] H. Li, D. Liu, D. Wang, and X. Yang, "Integral reinforcement learning for linear continuous-time zero-sum games with completely unknown dynamics," *IEEE Transactions on Automation Science and Engineering*, vol. 11, no. 3, pp. 706–714, 2014.
- [21] D. Vrabie and F. Lewis, "Integral reinforcement learning for online computation of feedback nash strategies of nonzero-sum differential games," in *Proceedings of IEEE Conference on Decision and Control* (CDC), Atlanta, GA, 2010.
- [22] R. Song, F. L. Lewis, and Q. Wei, "Off-policy integral reinforcement learning method to solve nonlinear continuous-time multiplayer nonzero-sum games," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 3, pp. 704–713, 2017.