

AUTONOMOUS AERIAL VEHICLES DISTRIBUTED CONTROL AND  
INTERACTIVE GAMES

by

YUSUF KARTAL

DISSERTATION

Submitted in partial fulfillment of the requirements  
for the degree of Doctor of Philosophy at  
The University of Texas at Arlington  
May, 2022

Arlington, Texas

Supervising Committee:

Frank Lewis, Supervising Professor

Atilla Dogan, Supervising Professor

Kamesh Subbarao

Animesh Chakravarthy

Shuo Linda Wang

Copyright by  
YUSUF KARTAL  
2022

*I dedicate this dissertation to my wife Tuğçe, my daughter, Gökçe, and my parents, Kamuran and Hami, for their constant support and enduring love.*

## ACKNOWLEDGEMENTS

Starting from my academic advisors Dr. Frank Lewis and Dr. Atilla Dogan, I would like to thank people who had an important role in my academic career. First of all, I would like to thank my supervisor, teacher, and friend, Dr. Frank Lewis, for inspiring my interest in the development of innovative technologies and encouraging me on carrying out on the edge research. Thanks to his expertise and our fruitful discussions, I have collected unforgettable memories in these years. Without his continuous support, it would have been very difficult for me to maintain my motivation through this dissertation. I would also like to thank Dr. Atilla Dogan for providing me opportunity to obtain Ph.D. in UTA.

I also want to thank the members of my dissertation committee, Dr. Kamesh Subbarao, Dr. Animesh Chakravarthy and Dr. Shuo Linda Wang, for their insightful comments that help to improve my research. Dr. Kamesh Subbarao and Dr. Animesh Chakravarthy deserve special thanks for teaching me important concepts in their classes that contribute towards finalizing this dissertation.

With my whole heart I thank all of my friends and collaborators at UTA. I would also like to extend my appreciation to Turkish Aerospace for providing me financial support. I am extremely grateful to Dr. Ugur Zengin for his unconditional support. I wish to thank the University of Texas at Arlington Research Institute family for being very helpful and kind to me.

Finally, I would like to express my deepest gratitude to my friends in Turkey, and my family. I am extremely fortunate to have them in my life. This work was supported in part by Army Research Office under Grant W911NF-20-1-0132, Office of Naval Research under Grant N00014-18-1-2221, and National Science Foundation under Grant 1839804, EAGER.

## LIST OF ILLUSTRATIONS

Figure	Page
1. Coordinate systems of the quadrotor. . . . .	12
2. Desired input states calculation for the attitude controller. . . . .	14
3. The flight controller FSM design. . . . .	31
4. The attitude control inputs of leader UAV . . . . .	33
5. Torque and thrust controls of the leader UAV . . . . .	34
6. The path tracked by the UAVs when there exists no delay and undirected graph topology is used . . . . .	35
7. The path tracked by the UAVs when there exists no delay and directed graph topology is used . . . . .	36
8. The path tracked by the UAVs when there exists two seconds delay and directed graph topology is used . . . . .	37
9. The path tracked by the UAV with the backstepping controller using 8-figure trajectory . . . . .	38
10. The path tracked by the UAV with the backstepping controller using a circular trajectory. . . . .	39
11. Controller behavior with time delay. . . . .	39
12. The path tracked by the UAVs with the distributed backstepping controller when the leader follows an eight-figure trajectory. . . . .	40
13. The path tracked by the UAVs with the distributed backstepping controller when the leader follows a circular trajectory. . . . .	41

14. Observation of time delay graph when the leader follows an eight-figure trajectory. . . . .	42
15. Observation of time delay graph when the leader follows a circular trajectory.	42
16. Sphere of collisions for players and their frames in 3-dimensions used for finite-time capture analysis. . . . .	59
17. Position of the pursuer and evader: (a) $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ , $\pi^e(\boldsymbol{\delta}) \triangleq \textit{suboptimal}$ , (b) $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ , $\pi^e(\boldsymbol{\delta}) \triangleq (3.18)$ . $L_2$ norm of (3.10): (c) $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ , $\pi^e(\boldsymbol{\delta}) \triangleq \textit{suboptimal}$ , (d) $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ , $\pi^e(\boldsymbol{\delta}) \triangleq (3.18)$ . . . . .	70
18. PE game when $\lambda^p = 5$ , $\lambda^e = 4$ : (a) game-optimal velocity of the pursuer and evader $\forall i \in \{x, y, z\}$ , (b) control force of the pursuer (3.47), (c) control force of the evader (3.47), (d) $L_2$ norm of (3.5) for each player, (e) Euler angles of the pursuer by (3.47), (f) Euler angles of the evader by (3.47), (g) body evaluated control force of the players (3.50), (h) weights $\forall i \in \{x, y, z\}$ convergence for the critic NN (3.45). . . . .	73
19. PE game when $\lambda^p = 10$ , $\lambda^e = 8$ : (a) game-optimal velocity of the pursuer and evader $\forall i \in \{x, y, z\}$ , (b) control force of the pursuer (3.47), (c) control force of the evader (3.47), (d) $L_2$ norm of (3.5) for each player, (e) Euler angles of the pursuer by (3.47), (f) Euler angles of the evader by (3.47), (g) body evaluated control force of the players (3.50), (h) weights $\forall i \in \{x, y, z\}$ convergence for the critic NN (3.45). . . . .	74
20. System response when no control is applied. . . . .	100
21. System response when the Nash gains (4.30) and (4.31) are employed. . .	101
22. Convergence of the gain matrix $\mathbf{K}_o^a$ parameters by using the Algorithms 1 and 2. . . . .	102
23. Convergence of the $\mathbf{N}_o$ matrix parameters by using the Algorithm 1. . . .	103
24. Convergence of the gain matrix $\mathbf{K}_o^a$ parameters by using the Algorithm 3.	104

25. Performance of the Algorithm 4. . . . .	105
26. Layout of Adversarial Leader-Follower Graphical Game . . . . .	110
27. The communication graph of the four followers and three leaders MLF game.	130
28. Path tracked by the UAVs in the MLF game. . . . .	132
29. Output containment position error for each follower in three dimensional space. . . . .	133
30. The Euler angles and total thrust of first follower UAV. . . . .	134
31. The Euler angles and total thrust of second follower UAV. . . . .	135
32. The Euler angles and total thrust of third follower UAV. . . . .	135
33. The Euler angles and total thrust of fourth follower UAV. . . . .	136



## LIST OF TABLES

Table	Page
1. Inner attitude control loop PID parameters . . . . .	32

## ABSTRACT

# AUTONOMOUS AERIAL VEHICLES DISTRIBUTED CONTROL AND INTERACTIVE GAMES

YUSUF KARTAL, Ph.D.

The University of Texas at Arlington, May, 2022

Supervising Professors: Frank Lewis, Atilla Dogan

As the number of quadrotors and other Unmanned Aerial Vehicles (UAVs) increases in industrial and urban areas, the development of reliable engineering methods to control their behavior as they interact with each other becomes of central interest in control research. With the increase in demand for UAVs to work together, several real-life challenges need to be addressed that include the implementation, test, and validation of the control algorithms in flight test experiments. On the one hand, the interests of UAVs may be in harmony to perform certain tasks such as transportation, surveillance and reconnaissance. On the other hand, the interests of UAVs may be directly opposite such as pursuing an evader UAV, and vice versa. This scenario is analyzed under category of the pursuit-evasion games in literature but constraints on the players' actions are not well considered. One of the primary requirements of autonomy is to be robust against the external disturbances, which can be achieved by designing a controller that guarantees  $L_2$  gain boundedness by a prescribed attenuation level. Unfortunately, the standard approaches to the such control problems result in the sub-optimal gain solutions for guaranteed stability. In addition, when multi-

agent leader-follower control is considered, the mutual interests among the followers can be addressed within a well-established game theoretic framework. In particular, this can be achieved via solving an output containment problem by introducing selfish followers where each follower only considers its own utility. However, standard approaches result in coupled Riccati equations that are hard to solve. Motivated by the desire to solve these problems, this dissertation has been written.

This dissertation first proposes a backstepping-based, distributed formation control method that is stable independent of time delays in communication among multiple UAVs. The proposed non-standard backstepping technique enables designer to develop an outer position & velocity control loop that interfaces seamlessly with the inner attitude controller of the cascaded control system for UAVs. Next, by using the nonstandard backstepping structure, we present a rigorous formulation for the pursuit-evasion (PE) game when velocity constraints are imposed on agents of the game or players. The game is formulated as an infinite-horizon problem using a non-quadratic functional, then sufficient conditions are derived to prove capture in a finite-time. Then, a new formulation for the  $H_\infty$  static output-feedback (OPFB) control problem that guarantees both stability and  $L_2$  gain boundedness of a Linear Time Invariant (LTI) system is given. This formulation allows us to extend multi-agent distributed formation control to multi-agent leader-follower (MLF) output containment game where the agents of game are UAVs. Lastly, the MLF output containment game is analyzed by introducing a novel cost functional whose solution provides both Nash and distributed robust control strategies in the sense that each follower uses the state information of its own and neighbors.

## TABLE OF CONTENTS

ACKNOWLEDGEMENTS . . . . .	iv
LIST OF ILLUSTRATIONS . . . . .	vi
LIST OF TABLES . . . . .	ix
ABSTRACT . . . . .	x
CHAPTER ONE: Introduction . . . . .	1
CHAPTER TWO: Distributed Backstepping Based Control of Multiple UAV Formation Flight Subject to Time Delays . . . . .	7
2.1. Preliminaries . . . . .	10
2.2. Backstepping control . . . . .	14
2.3. Distributed backstepping position control loop of multiple UAV with network delays . . . . .	22
2.4. Experiment design and flight test details . . . . .	30
2.5. Actual flight test results . . . . .	36
CHAPTER THREE: Optimal Game Theoretic Solution of the Pursuit-Evasion Intercept Problem Using On-Policy Reinforcement Learning . . . . .	44
3.1. Problem formulation and model description . . . . .	47
3.2. Developing velocity tracker using backstepping control method . . . . .	48
3.3. Optimal game theoretic velocity generation for pursuit-evasion game . . . . .	50
3.4. Online solution of HJI equation using integral reinforcement learning (IRL) . . . . .	61
3.5. Generalized rotational dynamics of the pursuer and evader . . . . .	65
3.6. Implementation on dynamic system . . . . .	69

CHAPTER FOUR: New Solution for H-infinity Static Output-Feedback Control	
Using Integral Reinforcement Learning . . . . .	75
4.1. Preliminaries and problem formulation . . . . .	78
4.2. New solution of OPFB $H_\infty$ game . . . . .	83
4.3. Online integral reinforcement learning solution algorithm for $H_\infty$ OPFB	92
4.4. Simulation results . . . . .	99
CHAPTER FIVE: Adversarial Multi-agent Output Containment Graphical Game	
with Local and Global Objectives for UAVs . . . . .	106
5.1. Preliminaries . . . . .	109
5.2. Multi-agent leader-follower game formulation . . . . .	117
5.3. Stability and $\mathcal{L}_2$ gain bound analysis with output feedback . . . . .	125
5.4. Simulation results . . . . .	130
CHAPTER SIX: Conclusions and Summary . . . . .	137
REFERENCES . . . . .	139
BIOGRAPHICAL STATEMENT . . . . .	153



## CHAPTER 1

### Introduction

Inspired by the naturally occurring biological groups such as herds and flocks where each member acts only under the influence of its neighbors [1], formation flight of quadrotors and other unmanned aerial vehicles (UAVs) has drawn great attention in recent years, due to their capability to perform certain tasks such as transportation [2], surveillance and reconnaissance, [3] and target search and detection [4]. With the increase in demand for UAVs to work together to accomplish these tasks, several real-life challenges need be addressed that include the implementation, test, and validation of the control algorithms in flight test experiments. In particular, we examine the challenge of designing a cooperative controller for UAVs to provide capabilities such as, performance despite communication time-delays in a leader-follower formation, employment of constrained input strategies, being robust against external disturbances, and adopting optimal actions to optimize mutual cost of the group.

In the second chapter of this dissertation, we propose a backstepping-based, distributed formation control method that is stable independent of time delays in communication among multiple UAVs. Instead of directly controlling the thrust generated by the propellers, we partition the mathematical model of the UAV into two subsystems, a linear attitude control loop and a nonlinear position control loop [5]. Centralized formation control of UAVs requires each agent to maintain a separation distance from other agents, which burdens the communication network of the UAVs. To overcome this problem, we consider a distributed control scheme wherein each agent updates its attitude and position based on the state information gathered

through its neighbors. Instead of directly controlling the thrust generated by the propellers, we partition the mathematical model of the UAV into two subsystems, a linear attitude control loop and a nonlinear position control loop. A backstepping-based outer position controller is then designed that interfaces seamlessly with the inner attitude controller of the cascaded control system. The closed-loop stability is established using a rigorous Lyapunov-Krasovskii analysis [6] under the influence of distributed network time delays. Using the directed graph topology and a distributed backstepping structure, it is shown that the stability criterion is delay-independent. The proposed control algorithms are verified in simulation and then implemented in hardware, and actual flight test experiments prove the validity of these algorithms.

This dissertation next copes with constrained input strategies for the pursuit-evasion (PE) game. A novel tracking Hamilton–Jacobi–Isaacs (HJI) equation associated with the non-quadratic value function is employed, which is solved for Nash equilibrium [7] velocity policies for each agent with arbitrary nonlinear dynamics. In contrast to the existing remedies for proof of capture in PE game, the proposed method does not assume players are moving with their maximum velocities and considers the velocity constraints a priori. Attaining the optimal actions requires the solution of HJI equations online and in real-time. We overcome this problem by presenting the on-policy iteration of Integral Reinforcement Learning (IRL) technique [8], [9]. The persistence of excitation for IRL to work is satisfied inherently until capture occurs, at which time the game ends. Furthermore, a nonlinear backstepping control method is proposed to track desired optimal velocity trajectories for players with generalized Newtonian dynamics. Simulation results are provided to show the validity of the proposed methods.

Then, a new formulation for the  $H_\infty$  [10]–[15] static output-feedback (OPFB) control problem that guarantees both stability and  $L_2$  gain boundedness of a Linear



Time Invariant (LTI) system is given. The problem is treated as a zero-sum differential game by introducing a quadratic performance index, and then a novel augmented Hamiltonian functional is proposed to solve for the Nash equilibrium point consisting of minimizing extrema (input) & maximizing extrema (disturbance) for the game of this kind. Unfortunately, the standard approaches to the  $H_\infty$  control problem with static OPFB result in the sub-optimal gain solutions for guaranteed stability [10]-[11], [16]. In this chapter, we provide necessary and sufficient conditions of the optimal gain solutions that inherently stabilize the system dynamics while also guaranteeing Nash equilibrium. To obtain the optimal gain solution, two off-line iterative solution algorithms are given. The first algorithm is based on Lyapunov iterations requires an initial stabilizing gain. A second algorithm based on Riccati iterations obviates the initial stabilizing gain requirement. Then, based on the Lyapunov iterations, an off-policy Integral Reinforcement Learning (IRL) algorithm [17]-[18] is developed to learn the optimal gain solution online without requiring any knowledge of system state, control, and disturbance matrices. Simulation results are provided to show the validity of the proposed methods.

Finally, multiple leader and follower graphical games that constitute challenging problems for aerospace and robotics applications are considered. One of the challenges is to address the mutual interests among the followers with an optimal control point of view. In particular, the traditional approaches [19]-[20], treat the output containment problem by introducing selfish followers where each follower only considers its own utility. In this chapter, we propose a differential output containment game over directed graphs where the mutual interests among the followers are addressed with an objective functional that also considers the neighboring agents. The obtained output containment error system results in a formulation where outputs of all followers are proved to converge to the convex hull spanned by the outputs of leaders [21] in

a game optimal manner. The output containment problem is solved using the  $\mathcal{H}_\infty$  output feedback method where the new necessary and sufficient conditions are presented. Another challenge is to design distributed Nash equilibrium control strategies for such games [22]-[26], which cannot be achieved with the traditional quadratic cost functional formulation. Furthermore, an  $\mathcal{L}_2$  gain bound of the output containment error system that experiences worst-case disturbances with respect to the  $\mathcal{H}_\infty$  criterion is investigated. The proposed methods are validated by means of multi-agent quadrotor Unmanned Aerial Vehicles (UAVs) output containment game simulations.

The resulting publications are listed below:

1. Yusuf Kartal , Kamesh Subbarao, Nicholas R. Gans, Atilla Dogan, Frank Lewis, Distributed backstepping based control of multiple UAV formation flight subject to time delays. *Published in IET Control Theory & Applications*. Volume 14, Issue 12. 2020.

<https://doi.org/10.1049/iet-cta.2019.1151>

2. Yusuf Kartal, Kamesh Subbarao, Atilla Dogan, Frank Lewis, Optimal game theoretic solution of the pursuit-evasion intercept problem using on-policy reinforcement learning. *Published in International Journal of Robust and Nonlinear Control*. Volume 31, Issue 16. 2021.

<https://doi.org/10.1002/rnc.5719>

3. Yusuf Kartal, Wenqian Xue, Atilla Dogan, Frank Lewis, New Solution for H-infinity Static Output-Feedback Control Using Integral Reinforcement Learning. 2021. *Under Review in IEEE Transactions on Cybernetics*.

4. Yusuf Kartal, Ahmet Taha Koru, Frank Lewis, Yan Wan, Atilla Dogan, Adversarial Multi-agent Output Containment Graphical Game with Local and Global Objectives for UAVs. 2021. *Under Review in IEEE Transactions on Control of Network Systems*.

The co-authors of above named work(s) acknowledge Yusuf Kartal as the primary author of the work(s) listed above. They authorize Yusuf Kartal to use the listed work(s) in this dissertation. They further agree that he may use this work to comply with requirements for graduation.

Distributed backstepping based control of multiple UAV formation flight subject to  
time delays [27]

Y. Kartal, K. Subbarao, N. Gans, A. Dogan and F. Lewis, *Published in IET Control  
Theory & Applications*. Volume 14, Issue 12. 2020.  
<https://doi.org/10.1049/iet-cta.2019.1151>

## CHAPTER 2

### Distributed Backstepping Based Control of Multiple UAV Formation Flight Subject to Time Delays

Formation control is a type of multi-agent architecture that relies on relative motion of agents [28]. There exist different approaches to ensure formation control for multi-agent systems in the control community. Three recognized categories are behavior-based approaches [29], virtual structure-based approaches [30] and leader-follower approaches [31]. In the behavior-based approaches, each agent of the formation acts according to predefined behavior. This approach is behaviorally inflexible since motion is predefined. Alternately, virtual structure-based approach introduces a virtual vehicle for each vehicle in the formation and transforms the formation problem into a trajectory tracking problem. However, since the virtual vehicles are not exposed to any type of disturbance in the environment, there is a high chance that the followers break formation in the event of unexpected environmental disturbances. On the other hand, the leader-follower approach is easy to implement, and all agents react to any environmental change, but the network delay must be examined well to maintain formation.

Several works deal with linear dynamics of multi-agent systems [32, 33, 34, 35, 36]. Particularly, [33] reveals some of the necessary and sufficient conditions to achieve predefined time-varying formations with switching interaction topologies based on the algebraic Riccati equation. The authors of [36] propose a distributed adaptive control technique that uses adaptive gain scheduling to tune the coupling weights between the individuals of the multi-agent system. Whereas, [37] uses the same technique to satisfy

prescribed  $H-\infty$  like performance and to manage the side effects of uncertainties in the system dynamics. However, since most physical systems are intrinsically nonlinear, these linear cooperative control methods cannot be applied directly [38]. Therefore, considerable number of works studied nonlinear dynamics of the multi-agent systems [38, 39, 40, 41, 42]. Neural-network based adaptive control and distributed impulsive control methods are examined to achieve leader-follower consensus with the class of nonlinear multi-agent systems in [38] and [39] respectively. The authors prove the stability of consensus error dynamics with well-known Lyapunov stability analysis techniques.

However, the problem of designing leader-follower formation strategies in which agents experience distributed network delays with pinning gain control requires more attention. One significant challenge associated with this relates to the construction of an appropriate compact nonlinear mathematical model of the multi-UAV system. In [6], authors study the effect of commensurate time delays on the closed-loop stability using the Retarded Functional Differential Equation (RFDE) form. By modeling the dynamics of the UAV as a double integrator, [43] and [44] treat each agent of the multi-UAV system as a point-mass system to apply time-varying consensus-based approaches on the multi-UAV system. This is a gross oversimplification of the UAV dynamics, especially for the quadrotor platforms considered in this paper. The inner/outer control loop partitioning allows us to deal with the delays occurring in the communication network of the UAVs, and to show that the stability of the outer control loop of the cascaded system is independent of delay. It is then shown that the closed-loop error dynamics of the whole system is also stable, independent of time delay.

Moreover, authors of [45, 46, 47, 48, 49] use nonlinear backstepping control method to deal with recursive design structure. This approach enables designer to

solve the stabilization problem partially for each submodule of the system of interest. Particularly, [45] proposes the adaptive backstepping control scheme to possess stronger stability properties while dealing with parametric uncertainties. [46] extends the adaptive backstepping control scheme in [45], to achieve both convergence to the path and predetermined dynamic behavior along the path simultaneously. In [47], integrator backstepping and quaternion feedback is adopted to stabilize the attitude of a micro satellite. Authors of [49] work on improving transient response of the closed-loop system by presenting a generalized backstepping process, based on the solvability of virtual controllers.

There are three major contributions of this work. Firstly, we use the second-order nonlinear dynamics of UAVs and synthesize a novel rigorous distributed backstepping control technique that has a form that easily extends to multiple UAV distributed control. Secondly, we partition the mathematical model of UAVs into two subsystems, an inner attitude control loop which is built into the quadrotor, and outer position controllers that consider relative motion of neighbors in formation flight. This allows us to rigorously analyze the delays occurring in the communication network of the UAVs. We prove that the stability of the outer control loop of the cascaded system is independent of delay, which implies that the closed-loop error dynamics of the whole system is also stable independent of delay. Lastly, we employ the directed graph topology to design formation control of the multi-UAV system, which reduces the work burden for the UAV communication network. It is rigorously shown that the stability criterion is delay-independent when each agent of the formation experiences distributed delays while communicating with its neighbors. The actual flight tests that show the validity and robustness of the developed control algorithms.

The rest of the paper is organized as follows. In Section 2.1, we provide the preliminaries of the mathematical model of the quadrotor and graph theory to under-

stand the basics of the distributed control approach. Section 2.2 brings an analysis of the control structures proposed, which involves the inner attitude controller and backstepping-based position controller for the task of trajectory tracking. We first explain the attitude controller design procedure for the quadrotors. Then we show the stability analysis of the backstepping control method. Section 2.3 illustrates how to extend the backstepping control algorithm to control multiple agents using the distributed backstepping tracker, which has a stable delay-independent system structure under the influence of non-constant distributed delays. Section 2.4 and 2.5 show the flight tests on a real UAV, where we illustrate the trajectories followed by AR.Drone 2.0 quadrotor with a full nonlinear backstepping tracker and by multiple AR.Drone 2.0's with a distributed backstepping trackers.

## 2.1 Preliminaries

The goal of this paper is to design a distributed controller for multi-UAV systems. Flying in a formation requires the agents to maintain separate distance from each other, which burdens the communication network and induces communication delays. In this section, we give preliminaries of a mathematical model of the quadrotor and graph theory to clarify the idea of the backstepping-based, distributed formation control method. In the next section, we present backstepping control for a single UAV. Then in Section 2.3, we present the formation controller with the network delays.

### 2.1.1 Mathematical model

This section introduces the standard nonlinear model of the quadrotor dynamics. To localize the quadrotor position, we use the Earth fixed frame. The origin of the three-dimensional (3D) axis system of the Body frame is assumed to be at



the center of mass of the quadrotor. The kinematics of the Euler angle rates can be expressed as

$$\mathbf{w}_B = \begin{bmatrix} p \\ q \\ r \end{bmatrix} = \begin{bmatrix} 1 & 0 & -s\theta \\ 0 & c\varphi & c\theta s\varphi \\ 0 & -s\varphi & c\theta c\varphi \end{bmatrix} \dot{\boldsymbol{\eta}} \quad (2.1)$$

where  $c$  and  $s$  refers to cosine and sine respectively and  $\mathbf{w}_B \in \mathbb{R}^3$  is the angular velocity in the Body frame components. Particularly,  $p$  is the roll rate,  $q$  is the pitch rate, and  $r$  is the yaw rate defined in the Body frame. Moreover,  $\boldsymbol{\eta} \in \mathbb{R}^3$  is the Euler angle vector (roll, pitch, and yaw) i.e.,  $\boldsymbol{\eta} = [\varphi \ \theta \ \psi]^T$ . Note that positive directions of Euler angles determined by right-hand rule, which are shown in Fig. 1.

The rotational dynamics are given by

$$\mathbf{I}_B \dot{\mathbf{w}}_B = \mathbf{S}(\mathbf{w}_B) \mathbf{I}_B \mathbf{w}_B + \boldsymbol{\tau}_B \quad (2.2)$$

where  $\mathbf{S}(\mathbf{w}_B) \in \mathbb{R}^{3 \times 3}$  is the skew-symmetric matrix [44],  $\boldsymbol{\tau}_B = [\tau_\varphi \ \tau_\theta \ \tau_\psi]^T$  is the torque vector and  $\mathbf{I}_B \in \mathbb{R}^{3 \times 3}$  is the inertia matrix defined in Body frame.

The translational dynamics of the quadrotor, ignoring any aerodynamic effects, is expressed in the Body frame is obtained to be

$$m \dot{\mathbf{U}} = \begin{bmatrix} 0 \\ 0 \\ \mu \end{bmatrix} + \mathbf{R} \mathbf{F}_g \quad (2.3)$$

where  $\mathbf{U} = [u \ v \ w]^T$  is the velocity vector defined in the Body frame,  $\mu$  is the total thrust produced by rotors in the Body frame  $z_B$ -axis. Note that  $m$  is mass of the

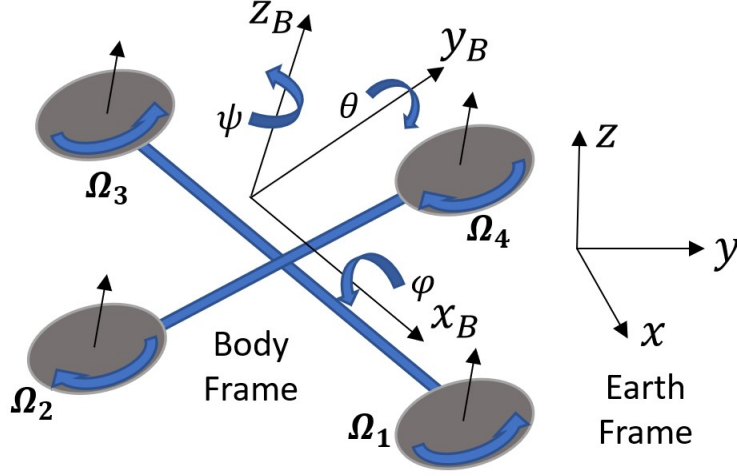


Figure 1: Coordinate systems of the quadrotor.

rigid body,  $\mathbf{F}_g = [0 \ 0 \ -mg]^T$  is the gravitational force vector and  $\mathbf{R} \in \mathbb{R}^{3 \times 3}$  is the rotation matrix from the Earth frame to the Body frame. Moreover  $u$ ,  $v$  and  $w$  stands for the velocities of the quadrotor in Body-axis coordinate system given in Fig.1. We obtain this rotation matrix using the yaw-pitch-roll (3-2-1) sequence. It is given by

$$\mathbf{R} = \begin{bmatrix} c\theta c\psi & c\theta s\psi & -s\theta \\ -c\phi s\psi + s\phi s\theta c\psi & c\phi c\psi + s\phi s\theta s\psi & s\phi c\theta \\ s\phi s\psi + c\phi s\theta c\psi & -s\phi c\psi + c\phi s\theta s\psi & c\phi c\theta \end{bmatrix} \quad (2.4)$$

Note that  $\mathbf{R}$  belongs to the special orthogonal group and is of rank 3, or  $SO(3)$ , whose determinant is equal to 1.

The translational dynamics of the quadrotor in the Earth frame is then formulated as

$$m\ddot{\boldsymbol{\xi}} = m\dot{\mathbf{V}} = \mathbf{F} + \mathbf{F}_g \quad (2.5)$$

where  $\boldsymbol{\xi} = [x \ y \ z]^T$  and  $\mathbf{V} \in \mathbb{R}^3$  denotes the position and velocity vectors in the Earth frame respectively. And,  $\mathbf{F} \in \mathbb{R}^3$  is the input force vector defined in the Earth frame.

Then 2.3 and 2.5 gives the following relation

$$\mathbf{F} = \begin{bmatrix} f_x \\ f_y \\ f_z \end{bmatrix} = \mathbf{R}^T \begin{bmatrix} 0 \\ 0 \\ \mu \end{bmatrix} = \begin{bmatrix} \mu(s\varphi s\psi + c\varphi s\theta c\psi) \\ \mu(-s\varphi c\psi + c\varphi s\theta s\psi) \\ \mu(c\varphi c\theta) \end{bmatrix} \quad (2.6)$$

### 2.1.2 Graph Theory

A Graph is constructed with a pair  $\mathbf{G} = (V, E)$ , where the set  $V = \{v_1, \dots, v_N\}$  defines the nodes or vertices, and  $E$  defines edges or arcs. The set  $E$  is composed of edge pairs  $(v_i, v_j)$ . If  $(v_i, v_j)$  is equal to  $(v_j, v_i) \forall i, j \in [0, N], i \neq j$ , then graph is said to be bidirectional. Each edge  $(v_j, v_i) \in E$ , has a weight  $a_{ij} > 0$  if and only if there exists a connection from node  $j$  to  $i$ . The graph is called undirected if  $a_{ij} = a_{ji}, \forall i, j$ . The undirected graph is said to be weight balanced, which leads to symmetric adjacency matrix  $\mathbf{A}$ .

The diagonal matrix  $\mathbf{D}$  is the  $i^{th}$  row sum of  $\mathbf{A}$  or weighted in-degree. Then, the Laplacian matrix is defined as

$$\mathcal{L} = \mathbf{D} - \mathbf{A}. \quad (2.7)$$

In this paper, the edge weights represent the trust between quadrotors, which are nodes of the formation graph. We create a graph topology based on adjacency or connectivity matrix  $\mathbf{A} = [a_{ij}]$ , realizing that  $a_{ii} = 0$ . The Laplacian matrices of all undirected graphs are real symmetric matrices. On the other hand, this is not valid for the digraphs. One of the contributions of this paper is proving consensus of the UAVs by adopting directed graph topology.

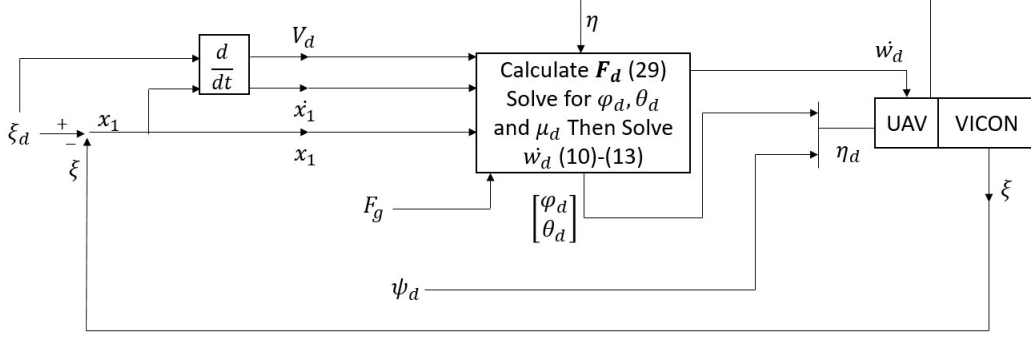


Figure 2: Desired input states calculation for the attitude controller.

## 2.2 Backstepping control

This section explains the full nonlinear backstepping control design for the quadrotor. The backstepping control structure derived here is shown in Fig.2. This is a nonstandard backstepping controller, has a novel form that allows direct extension to multiple interacting UAV control. Specifically, we derive the UAV error dynamics (2.33), which has a special form that is easily extended to multiple UAV formation control in Section 2.3. To apply the backstepping control method to the system defined in (2.5), we begin by adding and removing  $\mathbf{F}_d$ , an ideal virtual force input, and obtain the Newtonian model in terms of the desired forces

$$m\ddot{\boldsymbol{\xi}} = m\dot{\mathbf{V}} = \mathbf{F}_d + \mathbf{F}_g + \tilde{\mathbf{F}}_d \quad (2.8)$$

where  $\tilde{\mathbf{F}}_d = \mathbf{F} - \mathbf{F}_d$ . In Section 2.2.1, we show how to obtain desired Euler angle vector  $\boldsymbol{\eta}_d = [\varphi_d \ \theta_d \ \psi_d]^T$  to generate  $\mathbf{F}_d$ , and the time rate of change of desired vertical speed,  $\dot{w}_d$ . Then in Section 2.2.2,  $\boldsymbol{\tau}_B$  in (3.2) is designed using  $\boldsymbol{\eta}_d = [\varphi_d \ \theta_d \ \psi_d]^T$  and  $\dot{w}_d$  to get  $\tilde{\mathbf{F}}_d \rightarrow \mathbf{0}$ . Lastly in Section 2.2.3,  $\mathbf{F}_d$  is selected to get  $\boldsymbol{\xi} \rightarrow \boldsymbol{\xi}_d$ , where  $\boldsymbol{\xi}_d = [x_d \ y_d \ z_d]^T$  is the given desired position vector in the Earth frame. Proof of stability and tracking error convergence is given in Section 2.2.3.

### 2.2.1 Desired Euler angles

In Section 2.2.3, we show how to compute desired force vector  $\mathbf{F}_d$  to obtain position and velocity tracking. Herein we show how to compute  $\boldsymbol{\eta}_d$  and  $\dot{w}_d$  from the desired force data  $\mathbf{F}_d$  by using the inverse kinematics approach. Suppose we are given desired force  $\mathbf{F}_d$ . Note that from (2.6)

$$\mathbf{F}_d = \begin{bmatrix} f_{x_d} \\ f_{y_d} \\ f_{z_d} \end{bmatrix} = \begin{bmatrix} \mu_d(s\varphi_d s\psi_d + c\varphi_d s\theta_d c\psi_d) \\ \mu_d(-s\varphi_d c\psi_d + c\varphi_d s\theta_d s\psi_d) \\ \mu_d(c\varphi_d c\theta_d) \end{bmatrix} \quad (2.9)$$

where  $\mu_d$  is the desired thrust in the Body frame. (3.47) can be solved for the desired Euler angles

$$\tan\theta_d = \frac{f_{x_d}\cos\psi_d + f_{y_d}\sin\psi_d}{f_{z_d}},$$

$$\theta_d = \tan^{-1}\left(\frac{f_{x_d}\cos\psi_d + f_{y_d}\sin\psi_d}{f_{z_d}}\right), \quad (2.10)$$

$$\tan\varphi_d = \frac{\cos\theta_d (f_{x_d}\sin\psi_d - f_{y_d}\cos\psi_d)}{f_{z_d}},$$

$$\varphi_d = \tan^{-1}\left(\frac{\cos\theta_d (f_{x_d}\sin\psi_d - f_{y_d}\cos\psi_d)}{f_{z_d}}\right), \quad (2.11)$$

$$\mu_d = \frac{f_{z_d}}{\cos\varphi_d \cos\theta_d}. \quad (2.12)$$

and  $f_{z_d} \neq 0$ . Although it should be mentioned that  $f_{z_d} = 0$  only if  $\theta_d$  and/or  $\varphi_d = \pm\frac{\pi}{2}$  or  $\mu_d = 0$ . The condition  $\theta_d = \pm\frac{\pi}{2}$  or  $\varphi_d = \pm\frac{\pi}{2}$  correspond to singular orientations of the quadrotor and our domain of operation for  $\theta$  and  $\varphi$  is  $(-\frac{\pi}{2}, \frac{\pi}{2})$ . Further  $\mu_d = 0$  would correspond to zero total thrust. Furthermore, notice that  $\psi_d$  can be arbitrarily prescribed, and only the variables  $\theta_d$ ,  $\varphi_d$  and  $\mu_d$  must be found. The inner loop control design of the backstepping method requires the time rate of change of the

desired vertical speed in the Body frame,  $\dot{w}_d$ , which is calculated by using (2.3) such that

$$\dot{w}_d = \frac{\mu_d}{m} - g \cos \varphi_d \cos \theta_d . \quad (2.13)$$

Note that the information of time rate of change of the desired vertical speed or simply desired vertical acceleration acts as an input of the attitude controller loop will be given in Section 2.2.2. Derivation of this data is essential to control the height of the UAV, while accomplishing the path tracking objective accurately in 3D space.

### 2.2.2 Inner attitude control loop

In this section, we explain the inner attitude control of the backstepping method for the quadrotor. The attitude controller is generally built-in to the UAV and cannot be modified. This implies that the built-in attitude controller is assumed to track the quantities  $\boldsymbol{\eta}_d$  and  $\dot{w}_d$ .

We begin with deriving the desired Euler rates  $\boldsymbol{w}_{Bd} = [p_d \ q_d \ r_d]^T$  by using (2.1) and  $\boldsymbol{\eta}_d$  such that

$$\begin{bmatrix} p_d \\ q_d \\ r_d \end{bmatrix} = \begin{bmatrix} 1 & 0 & -s\theta_d \\ 0 & c\varphi_d & c\theta_d s\varphi_d \\ 0 & -s\varphi_d & c\theta_d c\varphi_d \end{bmatrix} \begin{bmatrix} \dot{\varphi}_d \\ \dot{\theta}_d \\ \dot{\psi}_d \end{bmatrix} . \quad (2.14)$$

Then the following PID controller is designed to generate changes in angular velocity of the propellers

$$\begin{bmatrix} \Delta\Omega_\varphi \\ \Delta\Omega_\theta \\ \Delta\Omega_\psi \end{bmatrix} = \begin{bmatrix} P_\varphi(\varphi_d - \varphi) + D_\varphi(p_d - p) + I_\varphi \int (\varphi_d - \varphi) \\ P_\theta(\theta_d - \theta) + D_\theta(q_d - q) + I_\theta \int (\theta_d - \theta) \\ P_\psi(\psi_d - \psi) + D_\psi(r_d - r) + I_\psi \int (\psi_d - \psi) \end{bmatrix}. \quad (2.15)$$

Furthermore, (3.59) is used to obtain the desired angular velocity of each rotor [2] such that

$$\begin{bmatrix} \Omega_{1d} \\ \Omega_{2d} \\ \Omega_{3d} \\ \Omega_{4d} \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 & 1 \\ 1 & -1 & 0 & -1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} \Omega_h + \Delta\Omega_{net} \\ \Delta\Omega_\varphi \\ \Delta\Omega_\theta \\ \Delta\Omega_\psi \end{bmatrix} \quad (2.16)$$

where  $\Omega_{id}$ ,  $i = 1, 2, 3, 4$ , corresponds to the desired angular velocities of the rotors and  $k$  is the thrust factor. And,  $\Omega_h$  is the rotor speed required to hover such that

$$\Omega_h = \sqrt{\frac{mg}{4k}} \quad (2.17)$$

$\Delta\Omega_{net}$  is the outcome of desired vertical acceleration in the Body frame,  $\dot{w}_d$  (2.13), in the form of

$$\Delta\Omega_{net} = \frac{m}{8k\Omega_h} \dot{w}_d. \quad (2.18)$$

Notice that while producing  $\Omega_h$  keeps the quadrotor at nominal condition (hover),  $\Delta\Omega_{net}$  moves the UAV along  $z_B$ -axis. In addition, producing  $\Delta\Omega_\varphi$ ,  $\Delta\Omega_\theta$  and  $\Delta\Omega_\psi$  deviates quadrotor from hover by resulting in roll, pitch and yaw respectively. Finally,

(3.60) yields the following torque vector expression, which stabilizes the rotational dynamics (3.2)

$$\boldsymbol{\tau}_B = \begin{bmatrix} \tau_\varphi \\ \tau_\theta \\ \tau_\psi \end{bmatrix} = \begin{bmatrix} lk(\Omega_{4d}^2 - \Omega_{2d}^2) \\ lk(\Omega_{3d}^2 - \Omega_{1d}^2) \\ d(\Omega_{1d}^2 + \Omega_{3d}^2 - \Omega_{2d}^2 - \Omega_{4d}^2) \end{bmatrix} \quad (2.19)$$

where  $l$  is the lever length, and  $d$  is the drag factor. The direction of angular velocities for each rotor is given in Fig.1, while first and third rotor turn anti-clockwise, the other two turn clockwise to cancel the yawing moments generated when the quadrotor is at nominal condition. The total thrust,  $\mu_d$ , is equal to the sum of thrusts generated by each rotor, that is

$$\mu_d = k(\Omega_{1d}^2 + \Omega_{2d}^2 + \Omega_{3d}^2 + \Omega_{4d}^2). \quad (2.20)$$

### 2.2.3 Outer position control loop

This section explains the outer position control loop of full nonlinear backstepping design for the quadrotor. We derive an expression for  $\mathbf{F}_d$  that guarantees the dynamics of (2.8) are stable. A main result is the form (2.33) for the error dynamics, which has a special structure that is directly extended to multiple quadrotors in Section 2.3.

Begin with defining the position and velocity errors in the Earth frame as

$$\begin{aligned} \mathbf{x}_1 &= \mathbf{e}_p = \boldsymbol{\xi}_d - \boldsymbol{\xi}, \\ \mathbf{x}_2 &= \mathbf{e}_v = \mathbf{V}_d - \mathbf{V} \end{aligned} \quad (2.21)$$

where  $\mathbf{V}_d = \dot{\boldsymbol{\xi}}_d \in \mathbb{R}^3$  is the desired velocity vector in the Earth frame.



The objective of the paper is to demonstrate the performance of the cooperative controller for any Commercial off the Shelf quadrotor platform with an inbuilt attitude controller. In such a situation, the assumption is reasonable that the inbuilt attitude controllers (typically PID) will accomplish this task. Thus the inner loop attitude controller is not analyzed further. Then to prove the convergence of error dynamics (3.45), which is the second step of backstepping method given in Section 2.2, we make following assumption

**Assumption 1.** *The inner attitude controller (3.43), tracks the Euler angles (3.48)-(3.49) and vertical acceleration (2.13). Hence the equilibrium of inner attitude control loop is stable.*

The next main theorem shows how to compute  $\mathbf{F}_d$  to guarantee stable position and velocity tracking of (2.8).

**Theorem 1.** *Under Assumption 1, the following control law, applied to the system governed by (2.8) ensures that the position and velocity tracking errors in (3.45)  $\rightarrow 0$  as  $t \rightarrow \infty$ .*

$$\begin{aligned} \mathbf{F}_d &= m\dot{\mathbf{V}}_d - \mathbf{F}_g + (m\mathbf{K}_1 + m\mathbf{K}_2) \dot{\mathbf{x}}_1 \\ &+ (m\mathbf{I}_{3 \times 3} + m\mathbf{K}_1\mathbf{K}_2) \mathbf{x}_1. \end{aligned} \quad (2.22)$$

*Proof.* The error dynamics are derived as

$$\dot{\mathbf{x}}_1 - \mathbf{x}_2 + \mathbf{x}_2^v - \mathbf{x}_2^v = \mathbf{0} \quad (2.23)$$

where  $\mathbf{x}_2^v$  is a virtual control signal. Moreover, by using (2.8) the error dynamics become

$$\dot{\mathbf{x}}_2 - \dot{\mathbf{V}}_d + \frac{\mathbf{F}_d}{m} + \frac{\mathbf{F}_g}{m} + \frac{\tilde{\mathbf{F}}_d}{m} = \mathbf{0}. \quad (2.24)$$

Then, define the following velocity error mismatch variable

$$\tilde{\mathbf{x}}_2 = \mathbf{x}_2^v - \mathbf{x}_2. \quad (2.25)$$

Substituting (2.25) in (2.23)

$$\dot{\mathbf{x}}_1 - \mathbf{x}_2^v = -\tilde{\mathbf{x}}_2. \quad (2.26)$$

Now we pick  $\mathbf{x}_2^v = -\mathbf{K}_1 \mathbf{x}_1$  where  $\mathbf{K}_1 \in \mathbb{R}^{3 \times 3}$  is a diagonal positive definite matrix. Then (2.26) becomes

$$\dot{\mathbf{x}}_1 + \mathbf{K}_1 \mathbf{x}_1 = -\tilde{\mathbf{x}}_2. \quad (2.27)$$

To examine the stability of (2.27), we pick the Lyapunov function candidate as follows

$$\mathcal{V} = \frac{1}{2} \mathbf{x}_1^T \mathbf{x}_1 + \frac{1}{2} \tilde{\mathbf{x}}_2^T \tilde{\mathbf{x}}_2. \quad (2.28)$$

Then the derivative of Lyapunov function candidate is derived using (2.24) and (2.26) as

$$\begin{aligned} \dot{\mathcal{V}} &= \mathbf{x}_1^T \dot{\mathbf{x}}_1 + \tilde{\mathbf{x}}_2^T \dot{\tilde{\mathbf{x}}}_2 \\ &= \mathbf{x}_1^T (-\mathbf{K}_1 \mathbf{x}_1 - \tilde{\mathbf{x}}_2) \end{aligned}$$

$$+ \widetilde{\mathbf{x}}_2^T \left( \dot{\mathbf{x}}_2^v - \dot{\mathbf{V}}_d + \frac{\mathbf{F}_d}{m} + \frac{\mathbf{F}_g}{m} + \frac{\widetilde{\mathbf{F}}_d}{m} \right). \quad (2.29)$$

To have strictly negative definite Lyapunov function derivative, we set  $\mathbf{F}_d$  as

$$\begin{aligned} \mathbf{F}_d &= m\dot{\mathbf{V}}_d - \mathbf{F}_g - m\dot{\mathbf{x}}_2^v + m\mathbf{x}_1 - m\mathbf{K}_2\widetilde{\mathbf{x}}_2 \\ &= m\dot{\mathbf{V}}_d - \mathbf{F}_g + (m\mathbf{K}_1 + m\mathbf{K}_2)\dot{\mathbf{x}}_1 \\ &\quad + (m\mathbf{I}_{3 \times 3} + m\mathbf{K}_1\mathbf{K}_2)\mathbf{x}_1 \end{aligned} \quad (2.30)$$

where  $\mathbf{K}_2 \in \mathbb{R}^{3 \times 3}$  is a diagonal positive definite matrix. Then, (2.29) becomes

$$\begin{aligned} \dot{\mathcal{V}} &= -\mathbf{x}_1^T \mathbf{K}_1 \mathbf{x}_1 - \widetilde{\mathbf{x}}_2^T \mathbf{K}_2 \widetilde{\mathbf{x}}_2 + \widetilde{\mathbf{x}}_2^T \frac{\widetilde{\mathbf{F}}_d}{m}, \\ &\leq -\lambda_{\min}(\mathbf{K}_1, \mathbf{K}_2) \|\tilde{\mathbf{x}}\|^2 + \widetilde{\mathbf{x}}_2^T \frac{\widetilde{\mathbf{F}}_d(0)}{m}, \\ &\leq -\lambda_{\min}(\mathbf{K}_1, \mathbf{K}_2) \|\tilde{\mathbf{x}}\|^2 + \|\widetilde{\mathbf{x}}_2\| \frac{\|\widetilde{\mathbf{F}}_d(0)\|}{m} \end{aligned} \quad (2.31)$$

where  $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x}_1^T & \widetilde{\mathbf{x}}_2^T \end{bmatrix}^T$  and  $\lambda_{\min}(\mathbf{K}_1, \mathbf{K}_2)$  stands for min eigenvalue of  $\mathbf{K}_1$  and  $\mathbf{K}_2$ . Note that in the worst case scenario,  $\lambda_{\min}(\mathbf{K}_1, \mathbf{K}_2)$  must be bigger than  $\frac{\|\widetilde{\mathbf{F}}_d(0)\|}{m}$ , which is a sufficient condition for asymptotic stability of origin. Moreover, from Assumption 1, the inner attitude control loop ensures that the quadrotor tracks the desired attitude angles  $\varphi_d$ ,  $\theta_d$ , and the desired thrust  $\mu_d$ , i.e  $\varphi \rightarrow \varphi_d$ ,  $\theta \rightarrow \theta_d$ , and  $\mu \rightarrow \mu_d$ . From the description of  $\mathbf{F}$  and  $\mathbf{F}_d$  in (2.6) and (3.47) respectively, we conclude that  $\mathbf{F} \rightarrow \mathbf{F}_d$  and hence,  $\widetilde{\mathbf{F}}_d \rightarrow \mathbf{0}$ . Then (2.31) becomes

$$\dot{\mathcal{V}} = -\mathbf{x}_1^T \mathbf{K}_1 \mathbf{x}_1 - \widetilde{\mathbf{x}}_2^T \mathbf{K}_2 \widetilde{\mathbf{x}}_2, \quad (2.32)$$

which is strictly negative definite since  $\mathbf{K}_1$  and  $\mathbf{K}_2$  are positive definite matrices. Note that  $\widetilde{\mathbf{x}}_2 \rightarrow \mathbf{0}$  implies  $\mathbf{x}_2 \rightarrow \mathbf{x}_2^v$ . Moreover,  $\mathbf{x}_2^v \rightarrow \mathbf{0}$  as  $\mathbf{x}_2^v = -\mathbf{K}_1 \mathbf{x}_1$ . Hence,  $\mathbf{x}_2 \rightarrow \mathbf{0}$  and the origin  $(\mathbf{0}, \mathbf{0})$ , which is the equilibrium of (3.45), is globally asymptotically stable.  $\square$

Using the control laws derived previously, from (2.30) the tracking error dynamics can be written in the state-space form as

$$\dot{\mathbf{x}} = \underbrace{\begin{bmatrix} \mathbf{0}_{3 \times 3} & \mathbf{I}_{3 \times 3} \\ -(\mathbf{K}_1 \mathbf{K}_2 + \mathbf{I}_{3 \times 3}) & -(\mathbf{K}_1 + \mathbf{K}_2) \end{bmatrix}}_{\mathbf{J}} \mathbf{x} \quad (2.33)$$

where  $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T \end{bmatrix}^T$ . Note that  $\mathbf{J}$  is Hurwitz. This form is instrumental in designing formation controllers for multiple UAV in the next section.

### 2.3 Distributed backstepping position control loop of multiple UAV with network delays

This section provides the connection of outer position control loop of backstepping method defined in Section 2.2.3 to the distributed multi-agent case. The error dynamics (2.33) are in a novel form which is easily extended in this section to multiple UAV formation.

We first treat the distributed system as delay-less. Then we perform rigorous stability analysis when agents experience both constant and distributed delays. Define dynamics (2.5) for each agent as

$$m\ddot{\boldsymbol{\xi}}_i = m\dot{\mathbf{V}}_i = \mathbf{F}_i + \mathbf{F}_g, \forall i = 1, \dots, N. \quad (2.34)$$

### 2.3.1 No communication delay

In this section, we first extend the error dynamics (2.33) to multiple quadrotors.

If there is no communication delay, define the position-based consensus error

$$\mathbf{e}_{\mathbf{p}i} = \sum_{j \in N_i} a_{ij} (\boldsymbol{\xi}_j - \boldsymbol{\Delta}_j - \boldsymbol{\xi}_i + \boldsymbol{\Delta}_i) + g_i (\boldsymbol{\xi}_0 - \boldsymbol{\xi}_i + \boldsymbol{\Delta}_i) \quad (2.35)$$

$\forall i = 1, \dots, N$  where  $g_i$  is the pinning gain,  $\boldsymbol{\Delta}_i$  (and  $\boldsymbol{\Delta}_j$ ) is the  $n$ -dim constant tracking offset vector of the  $i^{\text{th}}$  (and  $j^{\text{th}}$ ) UAV with respect to the  $n$ -dim position of the leader,  $\boldsymbol{\xi}_0 \in \mathbb{R}^3$ , of the formation. Lastly,  $N$  is the number of the UAVs in the formation. Note that  $g_i$  only takes values different than zero, if the node  $i$  is directly connected to the leader node. For the sake of simplicity, we use following vector notations

$$\begin{aligned} \mathbf{e}_p^c &= [\mathbf{e}_{p1}^T \ \mathbf{e}_{p2}^T \ \dots \ \mathbf{e}_{pN}^T]^T, \mathbf{e}_p^c \in \mathbb{R}^{Nn} \\ \underline{\boldsymbol{\Delta}} &= [\boldsymbol{\Delta}_1^T \ \boldsymbol{\Delta}_2^T \ \dots \ \boldsymbol{\Delta}_N^T]^T, \underline{\boldsymbol{\Delta}} \in \mathbb{R}^{Nn} \\ \boldsymbol{\xi}^c &= [\boldsymbol{\xi}_1^T \ \boldsymbol{\xi}_2^T \ \dots \ \boldsymbol{\xi}_N^T]^T, \boldsymbol{\xi}^c \in \mathbb{R}^{Nn} \\ \mathbf{V}^c &= [\mathbf{V}_1^T \ \mathbf{V}_2^T \ \dots \ \mathbf{V}_N^T]^T, \mathbf{V}^c \in \mathbb{R}^{Nn} \\ \mathbf{F}^c &= [\mathbf{F}_1^T \ \mathbf{F}_2^T \ \dots \ \mathbf{F}_N^T]^T, \mathbf{F}^c \in \mathbb{R}^{Nn}. \end{aligned} \quad (2.36)$$

By using (2.5), the global system dynamics for followers can be written as

$$m\dot{\mathbf{V}}^c = \mathbf{F}^c + \mathbf{1}_N \otimes \mathbf{F}_g. \quad (2.37)$$

Then, by using (2.7) and noting the fact  $\mathcal{L}\mathbf{1}_N = \mathbf{0}_N$  since the row sum of  $\mathcal{L}$  is zero, re-write (2.35) as

$$\mathbf{e}_p^c = -((\mathcal{L} + \mathbf{G}) \otimes \mathbf{I}_n) (\boldsymbol{\xi}^c - \underline{\boldsymbol{\Delta}}) + (\mathbf{G} \otimes \mathbf{I}_n)(\mathbf{1}_N \otimes \boldsymbol{\xi}_0)$$

$$= -((\mathcal{L} + \mathbf{G}) \otimes \mathbf{I}_n) (\boldsymbol{\xi}^c - \underline{\boldsymbol{\Delta}} - \mathbf{1}_N \otimes \boldsymbol{\xi}_0) \quad (2.38)$$

where  $\mathbf{1}_N$  is  $N$ -dim vector whose all elements are ones and  $\mathbf{G} \in \mathbb{R}^{N \times N}$  is a diagonal pinning gain matrix with the diagonal elements of  $g_i \forall i = 1, \dots, N$ . In addition,  $\otimes$  stands for the Kronecker product. Define  $\underline{\boldsymbol{\xi}}_0 = \mathbf{1}_N \otimes \boldsymbol{\xi}_0$ . Then (2.38) becomes

$$\mathbf{e}_p^c = ((\mathcal{L} + \mathbf{G}) \otimes \mathbf{I}_n) (\underline{\boldsymbol{\xi}}_0 + \underline{\boldsymbol{\Delta}} - \boldsymbol{\xi}^c). \quad (2.39)$$

The velocity-based consensus error is

$$\mathbf{e}_v^c = ((\mathcal{L} + \mathbf{G}) \otimes \mathbf{I}_n) (\underline{\mathbf{V}}_0 - \mathbf{V}^c) \quad (2.40)$$

where  $\underline{\mathbf{V}}_0 = \dot{\underline{\boldsymbol{\xi}}}_0 \in \mathbb{R}^{Nn}$ , then the error dynamics are derived as

$$\begin{aligned} \dot{\mathbf{e}}_p^c &= \mathbf{e}_v^c, \\ m\dot{\mathbf{e}}_v^c &= ((\mathcal{L} + \mathbf{G}) \otimes \mathbf{I}_n) (m\underline{\dot{\mathbf{V}}}_0 - m\dot{\mathbf{V}}^c). \end{aligned} \quad (2.41)$$

Substituting (2.37) in (2.41) results in

$$m\dot{\mathbf{e}}_v^c = ((\mathcal{L} + \mathbf{G}) \otimes \mathbf{I}_n) (m\underline{\dot{\mathbf{V}}}_0 - (\mathbf{F}^c + \mathbf{1}_N \otimes \mathbf{F}_g)). \quad (2.42)$$

Now, we set the global desired force vector that contains desired force information for each agent of the formation (2.36) by using mixed-product property of Kronecker product,  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$  such that

$$\begin{aligned} \mathbf{F}^c &= m\underline{\dot{\mathbf{V}}}_0 - \mathbf{1}_N \otimes \mathbf{F}_g + m(\mathbf{I}_N \otimes (\mathbf{K}_1 + \mathbf{K}_2))\mathbf{e}_v^c \\ &\quad + m(\mathbf{I}_N \otimes (\mathbf{I}_{3 \times 3} + \mathbf{K}_1\mathbf{K}_2))\mathbf{e}_p^c. \end{aligned} \quad (2.43)$$

Note that  $\mathbf{F}^c$  is the global form of (2.30). Moreover,  $\underline{\boldsymbol{\xi}}_0 + \underline{\boldsymbol{\Delta}}$  and  $\underline{\mathbf{V}}_0$  are the global form of  $\boldsymbol{\xi}_d$  and  $\mathbf{V}_d$  respectively. Then, we end up with the following second-order error dynamics such that

$$\begin{aligned}\dot{\mathbf{e}}_p^c &= \mathbf{e}_v^c \\ \dot{\mathbf{e}}_v^c &= -((\mathcal{L} + \mathbf{G}) \otimes (\mathbf{K}_1 \mathbf{K}_2 + \mathbf{I}_{3 \times 3})) \mathbf{e}_p^c \\ &\quad - ((\mathcal{L} + \mathbf{G}) \otimes (\mathbf{K}_1 + \mathbf{K}_2)) \mathbf{e}_v^c.\end{aligned}\tag{2.44}$$

The state-space form of (2.44) is

$$\dot{\mathbf{x}}^c(t) = \mathbf{J}^c \mathbf{x}^c(t)\tag{2.45}$$

where  $\mathbf{J}^c \in \mathbb{R}^{2Nn \times 2Nn}$  is the global system matrix such that

$$\mathbf{J}^c = \begin{bmatrix} \mathbf{0}_{Nn \times Nn} & \mathbf{I}_{Nn \times Nn} \\ -((\mathcal{L} + \mathbf{G}) \otimes (\mathbf{K}_1 \mathbf{K}_2 + \mathbf{I}_{3 \times 3})) & -((\mathcal{L} + \mathbf{G}) \otimes (\mathbf{K}_1 + \mathbf{K}_2)) \end{bmatrix}\tag{2.46}$$

and  $\mathbf{x}^c(t) = \begin{bmatrix} \mathbf{e}_p^{cT} & \mathbf{e}_v^{cT} \end{bmatrix}^T \in \mathbb{R}^{2Nn \times 1}$  is the global state vector. The global dynamics (2.45) are the combination of the single-agent dynamics (2.33) for the entire formation. Before we do the stability analysis for the closed-loop error dynamics given in (2.45), we make following assumption:

**Assumption 2.** *The graph topology of the multi-agent system contains a spanning tree with the root node being the leader node. This means that there is a directed path (not necessarily unique) from the leader node to every follower node.*

The next theorem extends the single-agent result in Theorem 1 to the multi-agent case by using M-matrix properties of the digraphs [50].

**Theorem 2.** *Given the Assumption 2,  $\mathcal{L} + \mathbf{G}$  is an irreducible M-matrix and has all eigenvalues strictly in the open right-half plane [50]. Then, the equilibrium of closed-*

loop error dynamics given in (2.45) is globally asymptotically stable point meaning that  $\mathbf{J}^c$  is Hurwitz.

*Proof.* Use the fact that Kronecker product of a positive diagonal matrix and an M-matrix has all eigenvalues strictly in the open right-half plane [51, 52]. Then, re-write (2.45) in form of the second-order differential equation such that

$$\begin{aligned} \ddot{\mathbf{e}}_p^c + ((\mathcal{L} + \mathbf{G}) \otimes (\mathbf{K}_1 + \mathbf{K}_2)) \dot{\mathbf{e}}_p^c \\ + ((\mathcal{L} + \mathbf{G}) \otimes (\mathbf{K}_1 \mathbf{K}_2 + \mathbf{I}_{3 \times 3})) \mathbf{e}_p^c = \mathbf{0}. \end{aligned} \quad (2.47)$$

Notice that all coefficient matrices of the characteristic polynomial of (2.47) have eigenvalues at open right half plane, hence the origin is globally asymptotically stable equilibrium by Routh-Hurwitz test. Note that if the graph topology was undirected, algorithms proposed in this paper would still work because  $\mathcal{L} + \mathbf{G}$  would be positive definite symmetric matrix with the assumption of there exists a path from the leader node to every follower node.  $\square$

### 2.3.2 Communication delays

In this section, first we consider the system with constant communication delay, which occurs while the local positioning system shares the position of each agent to their neighbors. This delay may be created by processing time of the positioning system, data header analysis, and storage at routers, etc. Note that this delay is upper bounded by the practical limitations. Then (2.45) is written in the form of Retarded Functional Differential Equation [6] such that

$$\dot{\mathbf{x}}^c(t) = \mathbf{J}^c \mathbf{x}^c(t) + \beta \mathbf{J}^c \mathbf{x}^c(t - \gamma) \quad (2.48)$$



where  $\gamma$  is the network delay,  $\beta$  is the gain of delayed term and  $\mathbf{J}^c$  is the system matrix (2.46).

**Theorem 3.** *For the system in (2.48), the origin is stable equilibrium for  $\beta \in (-1, 1]$  as the system matrix  $\mathbf{J}^c$  is Hurwitz by Theorem 2.*

*Proof.* As  $\gamma \rightarrow 0$ ,  $\beta$  must be greater than -1 so that the overall system is stable, which is the lower bound of  $\beta$ .

To find the upper bound, use the fact that  $\rho((jw\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) < 1 \forall w > 0$  as given in [53] where  $\rho(\cdot)$  denotes the spectral radius of a matrix and  $w$  denotes the frequency.

First assume that  $\rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) = 1, \forall w_t > 0$ , which implies  $e^{j\sigma_t}$  is the eigenvalue of the matrix  $(jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c$  for  $\sigma_t \in [0, 2\pi]$ . Then,  $\det(\mathbf{I} - (jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c e^{j\gamma_t w_t}) = 0$  for  $\gamma_t = \frac{\sigma_t}{w_t}$ , or equivalently by using matrix determinant lemma

$$\det(jw_t\mathbf{I} - \mathbf{J}^c - \mathbf{J}^c e^{j\gamma_t w_t}) = 0. \quad (2.49)$$

Hence (2.48) is not stable independent of delay with this assumption.

Next, assume that  $\rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) > 1, \forall w_t > 0$ . Since  $\rho(\cdot)$  is continuous function of  $w_t$  and

$$\lim_{w_t \rightarrow \infty} \rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) = 0. \quad (2.50)$$

Then  $\exists w_t \in (w, \infty)$  such that  $\rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) = 1$ , which makes (2.48) not stable independent of delay as this ends up with (2.49). Consequently, we show that  $\rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) < 1, \forall w_t > 0$ . By using Gelfand Corollaries, this results in

$$\begin{aligned} \rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1}\mathbf{J}^c) &\leq \rho((jw_t\mathbf{I} - \mathbf{J}^c)^{-1})\rho(\mathbf{J}^c) \\ &\leq \frac{\|\mathbf{J}^c\|_\infty}{\sqrt{\|\mathbf{J}^c\|_\infty^2 + w_t^2}}. \end{aligned} \quad (2.51)$$

Then the upper bound of  $\beta$  must be 1 by (2.51) to have  $\rho(\beta(j\omega_t \mathbf{I} - \mathbf{J}^c)^{-1} \mathbf{J}^c) < 1$ . To this end, we proved that as  $\beta \in (-1, 1]$ , the system represented in (2.48) is stable independent of delay.  $\square$

Now, if we consider the system with non-constant distributed delays, (2.45) can be written as

$$\dot{\mathbf{x}}^c(t) = \mathbf{J}^c \mathbf{x}^c(t) + \beta \int_{-\gamma}^0 \mathbf{J}^c \mathbf{x}^c(t+s) ds \quad (2.52)$$

where  $s \in [-\gamma, 0]$ ,  $\beta$  is the gain of delayed term and  $\gamma$  is the maximum delay.

**Theorem 4.** *The origin is an asymptotically stable equilibrium of (2.52) when there exist distributed delays in the communication network with gain  $\beta \in [0, 1]$ , which is a sufficient condition for asymptotic stability.*

*Proof.* As shown in Theorem 3, the origin is stable equilibrium for  $\beta \in (-1, 1]$  since  $\mathbf{J}^c$  is proven to be Hurwitz in Theorem 2. With this in mind, pick Lyapunov-Krasovskii functional as

$$\begin{aligned} V(\mathbf{x}_t^c) &= \mathbf{x}^{cT}(t) \mathbf{P} \mathbf{x}^c(t) \\ &+ \beta \int_{t-\gamma}^t \left[ \int_s^0 \mathbf{x}^{cT}(\ell) \mathbf{S} \mathbf{x}^c(\ell) d\ell \right] ds \end{aligned} \quad (2.53)$$

where  $\mathbf{P} \in \mathbb{R}^{2Nn \times 2Nn}$  and  $\mathbf{S} \in \mathbb{R}^{2Nn \times 2Nn}$  are positive definite, symmetric matrices. To have strictly positive definite Lyapunov-Krasovskii functional (2.53), the sufficient condition is  $\beta > 0$ . Before taking the derivative of Lyapunov-Krasovskii functional and developing stability analysis, we use change of variable  $\mathbf{f}(T) = \mathbf{x}^c(t+T)$  for arbitrary  $T$ , to simplify the stability analysis. Then, (2.52) and (2.53) become

$$\dot{\mathbf{f}}(0) = \mathbf{J}^c \mathbf{f}(0) + \beta \int_{-\gamma}^0 \mathbf{J}^c \mathbf{f}(s) ds \quad (2.54)$$

$$\begin{aligned}
\mathbf{V}(\mathbf{f}) &= \mathbf{f}^T(0) \mathbf{P} \mathbf{f}(0) \\
&+ \beta \int_{-\gamma}^0 \left[ \int_s^0 \mathbf{f}^T(\ell) \mathbf{S} \mathbf{f}(\ell) d\ell \right] ds.
\end{aligned} \tag{2.55}$$

Furthermore, by using the Leibniz Integral Rule and (2.54), the derivative of Lyapunov-Krasovskii functional (2.55) becomes

$$\begin{aligned}
\dot{\mathbf{V}}(\mathbf{f}) &= \mathbf{f}^T(0) [\mathbf{P} \mathbf{J}^c + \mathbf{J}^{cT} \mathbf{P} + \gamma \beta \mathbf{S}] \mathbf{f}(0) \\
&+ 2 \mathbf{f}^T(0) \int_{-\gamma}^0 \mathbf{P} \mathbf{J}^c \mathbf{f}(s) ds \\
&- \beta \int_{-\gamma}^0 \mathbf{f}^T(s) \mathbf{S} \mathbf{f}(s) ds.
\end{aligned} \tag{2.56}$$

To facilitate further development, (2.56) is written as

$$\begin{aligned}
\dot{\mathbf{V}}(\mathbf{f}) &= \mathbf{f}^T(0) [\mathbf{P} \mathbf{J}^c + \mathbf{J}^{cT} \mathbf{P}] \mathbf{f}(0) \\
&+ \int_{-\gamma}^0 \mathbf{v}^T \begin{bmatrix} \beta \mathbf{S} & \mathbf{P} \mathbf{J}^c \\ \mathbf{J}^{cT} \mathbf{P} & -\beta \mathbf{S} \end{bmatrix} \mathbf{v} ds
\end{aligned} \tag{2.57}$$

where  $\mathbf{v} = \begin{bmatrix} \mathbf{f}(0)^T & \mathbf{f}(s)^T \end{bmatrix}^T$ . Notice that  $\mathbf{V}(\mathbf{f}) \geq \varepsilon \left\| \mathbf{f}(0) \right\|^2$  is satisfied for sufficiently small  $\varepsilon > 0$ . And,  $[\mathbf{P} \mathbf{J}^c + \mathbf{J}^{cT} \mathbf{P}] \leq -\varepsilon \mathbf{I}$  since  $\mathbf{J}^c$  is proved to be Hurwitz by Theorem 2. In addition, assuming  $\exists \mathbf{P} = \mathbf{P}^T > \mathbf{0}$  and using the linear matrix inequality [13], the negative definiteness of a matrix  $\begin{bmatrix} \beta \mathbf{S} & \mathbf{P} \mathbf{J}^c \\ \mathbf{J}^{cT} \mathbf{P} & -\beta \mathbf{S} \end{bmatrix}$  implies

$\dot{\mathbf{V}}(\mathbf{f}) \leq \varepsilon \left\| \mathbf{f}(0) \right\|^2$ . Therefore, all conditions of the asymptotically stability by analyzing the derivative of Lyapunov-Krasovskii functional given in [23], have met meaning that the origin is asymptotically stable equilibrium.  $\square$

## 2.4 Experiment design and flight test details

This section addresses the crucial elements of our experiments, which are lab environment, flight controller design and simulations. An actual flight test is conducted in Section 2.5.

### 2.4.1 Lab environment

Equipment used are the Vicon, Parrot AR.Drone 2.0, and the master computer. Vicon is a motion capture system that provides the position of the UAVs.

The communication between master computer and Vicon is done via User Datagram Protocol (UDP). The frequency of the UDP Packets taken from the Vicon motion capture system is 100 Hz. The AR.Drone 2.0 has a built-in gyroscope and Inertial Measurement Unit (IMU) sensor suite. In practical applications, many quadrotors are designed with a built-in attitude controller and AR.Drone has its own attitude controller. This controller takes the desired values of  $\varphi_d$ ,  $\theta_d$  and  $\psi_d$  as inputs. The communication between the master computer and the AR.Drone is done via UDP. The frequency of UDP packages is set to 500 Hz. MATLAB-Simulink is used to create UDP nodes that are communicating with AR.Drone and Vicon. The receiver and the sender UDP nodes are inserted to the Simulink model in the form of S-functions. The controller and the trajectory generation algorithms are implemented in the model. Simulink-Desktop Real Time Add-on is used to send the real time commands to the quadrotor. The UDP nodes tolerate up to 10% packet loss rate, which is necessary to handle communication channel noise created by the lab environment.

### 2.4.2 Flight controller design

The flight controller is a high-level decision-making mechanism that activates different modes of operation depending on the state of the UAV. We recognize three

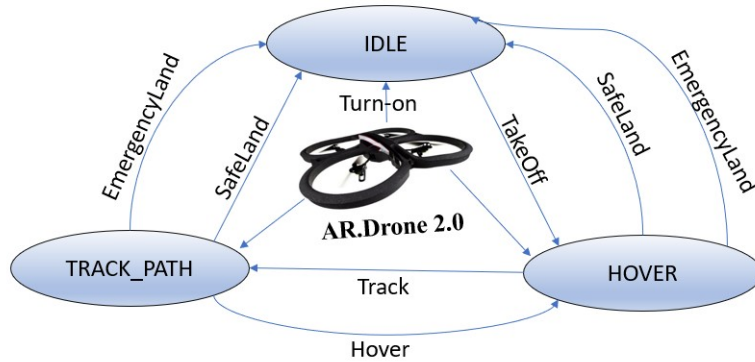


Figure 3: The flight controller FSM design.

modes of operation in our MATLAB implementation, which are IDLE, HOVER and TRACK\_PATH as shown in Fig. 3. The quadrotor enters the IDLE mode, when either we turn-on the AR.Drone manually or it is landed by receiving the Land command. While in the IDLE, the UAV is actively receiving the data packets via UDP port communication and is ready to get the TakeOff command.

When the UAV reaches the desired height,  $z_d$ , the flight controller switches to the HOVER mode. In this mode the UAV is at the nominal condition, its attitude is parallel to ground and motionless in the air. If the Track command is received in the HOVER mode, the flight controller switches to the PATH.TRACK mode. This mode is triggered after 10 seconds passed from the transmission of TakeOff command. In this mode, AR.Drone begins to track the predetermined trajectory. The flight control algorithm, first reads the IMU sensor buffers and then Vicon buffer to construct the close loop error dynamics. To calculate desired Euler angles, (3.48)-(3.50) is used. Notice herein the desired yaw angle command is set to an arbitrary constant.

To land the quadrotor on the ground, we use either EmergencyLand or SafeLand commands. The difference of these two commands is the timing of stopping propeller movements. If we send the SafeLand command to the quadrotor, it reduces the

Gains	$\Phi(Roll)$	$\theta(Pitch)$	$\psi(Yaw)$
<b>P</b>	6.42	6.42	4.82
<b>D</b>	5.54	5.54	7.89
<b>I</b>	1.85	1.85	0.11

Table 1: Inner attitude control loop PID parameters

propellers' speed till the height is in the range of 0-0.1 meter and shuts down the propellers. Else if we send the EmergencyLand to the quadrotor, it directly stops the propellers and lands on the ground. The appropriate structure for implementing the flight controller is the finite state machine (FSM) since the mode switching event is driven as shown in Fig. 3.

### 2.4.3 Simulations

The aim of this section is to verify control algorithms proposed in this paper by conducting different test scenarios. Before we implement the distributed backstepping control algorithm in the actual hardware, mathematical model in Section 2.1.1, Backstepping control method in Section 2.2 and distributed backstepping trackers in Section 2.3 are implemented in the Simulink. We first verified inner attitude control loop design in Section 2.2.2, the PID gains are given in Table 1.

To derive PID gains given in Table 1, thrust, drag factors, mass and arm length of the AR.Drone 2.0 must be measured. Mass and arm length of the quadrotor are measured as  $m = 0.467kg$  and  $l = 0.1785m$  respectively. To determine the thrust factor  $k$ , we first measure the angular velocity of a propeller with tachometer when the quadrotor is in hover. Then, we used (3.41) to derive  $k$ , which is found as  $8 * 10^{-6} N * s^2/rad^2$ . After that, using the thrust ratio analysis for small UAVs, the thrust factor  $d$  is derived as  $2 * 10^{-7} N * m * s^2/rad^2$ . Moreover the diagonal elements of  $\mathbf{K}_1$ ,  $\mathbf{K}_2$  are tuned as 2, 2, 3 and 1.5, 1.5, 3 respectively.

To construct the desired circular trajectory for the formation leader,  $x_d$  is set to  $\cos(\omega_t(t - t_{track}))$  and  $y_d$  is set to  $\sin(\omega_t(t - t_{track}))$  where  $t$  is simulation time,  $t_{track}$  is the time at which the UAV begins to track circular trajectory, and  $\omega_t$  stands for the frequency of the sinusoidal function. In our simulations, we pick  $t_{track}$  as 15s and  $\omega_t$  as  $0.5rad/s$ . Note that before formation leader begins to track circular trajectory, the  $x_d$  value is linearly increased by  $1m$  for  $t \in [10, 15]$ . Therefore, the leader UAV begins to track circular trajectory at 15s. In Section 2.5.1, we double value of  $\omega_t$  for  $y_d$  setting to construct the eight figure trajectory.

#### 2.4.3.1 No delay, undirected graph

The aim of this section is to show validity of the proposed algorithms, when the undirected graph topology is adopted to design leader-follower formation control while multi-UAV system does not experience any delays. We share the control histories for the formation leader in Fig. 4 and Fig. 5. Particularly, in Fig. 4, we show attitude

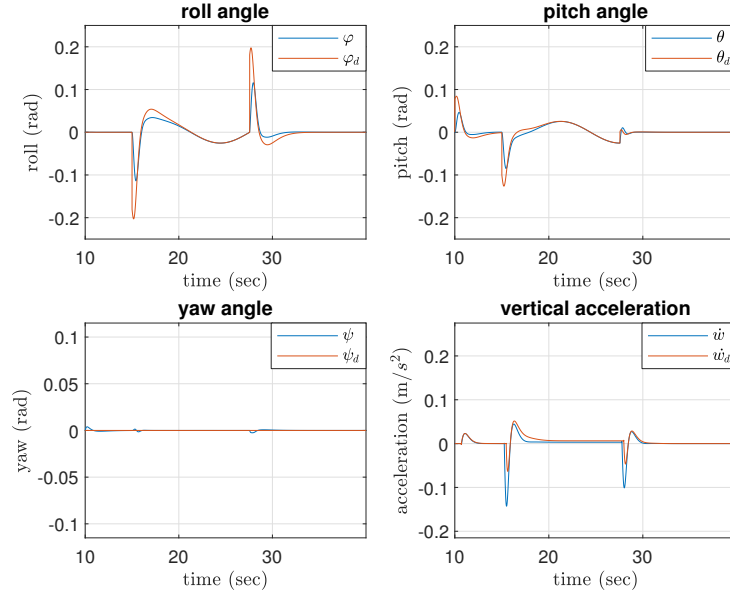


Figure 4: The attitude control inputs of leader UAV

control inputs (3.48)-(2.13) of the leader UAV. In Fig. 5, we show torque (3.43), and thrust (3.44) controls defined in the Body frame of the formation leader. For this test scenario, we pick Adjacency, pinning gain matrices and offset vector as follows

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \mathbf{G} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

$$\underline{\Delta}^T = \begin{bmatrix} 1.5 & 0 & 0 & 3 & 0 & 0 & 4.5 & 0 & 0 \end{bmatrix}. \quad (2.58)$$

Fig. 6 shows the leader and follower positions with the Adjacency and pinning gain matrices given in (2.58), when there is no communication delay in the multi-UAV communication network.

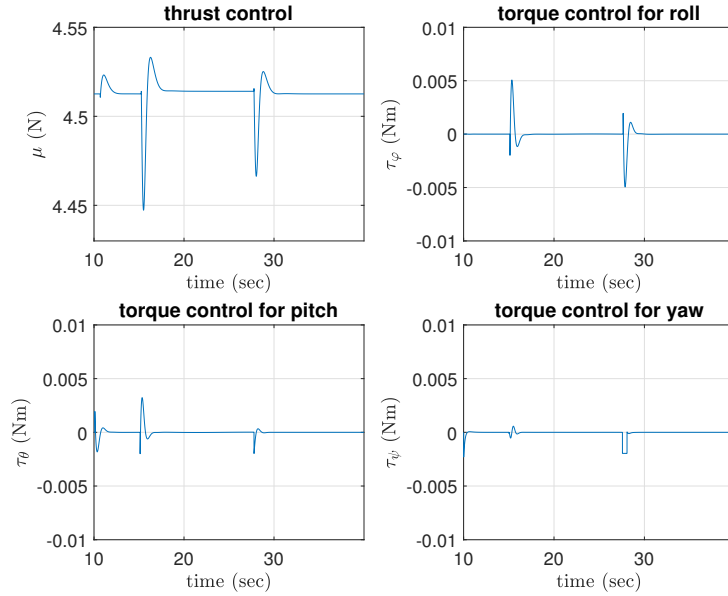


Figure 5: Torque and thrust controls of the leader UAV



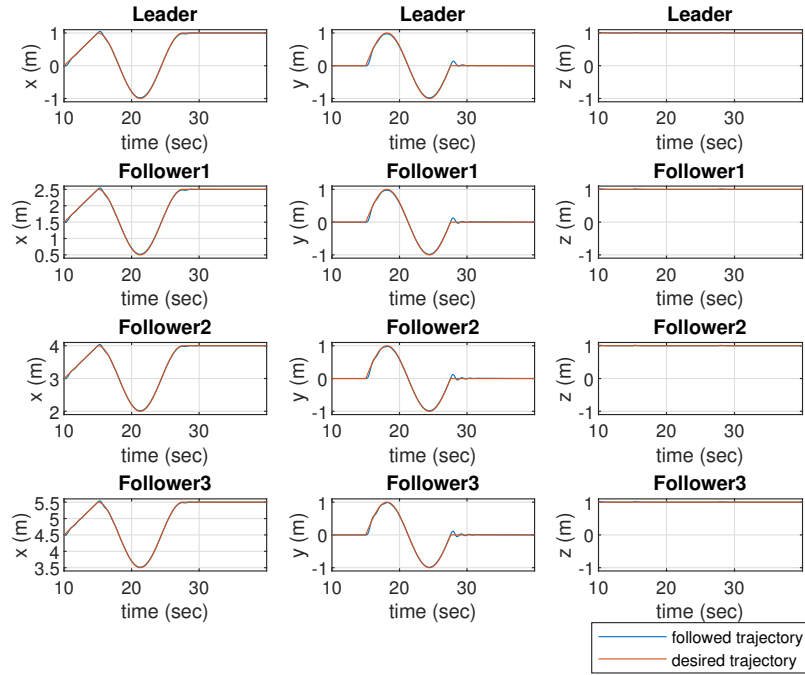


Figure 6: The path tracked by the UAVs when there exists no delay and undirected graph topology is used

#### 2.4.3.2 No delay, directed graph

For this test scenario, along with the offset vector given in (2.58) we pick Adjacency and pinning gain matrices that do not contradict Assumption 2, as follows

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \mathbf{G} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}. \quad (2.59)$$

Note that Fig. 7 is the same as Fig. 6 since the graph topology with the Adjacency and pinning gain matrices given in (5.56), contains a spanning tree.

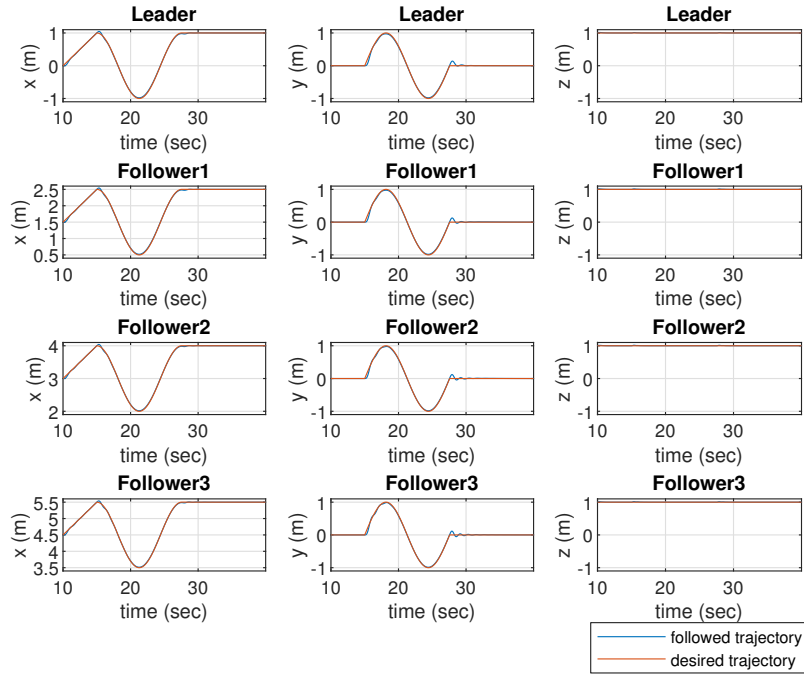


Figure 7: The path tracked by the UAVs when there exists no delay and directed graph topology is used

#### 2.4.3.3 With delay, directed graph

For the last test scenario of simulations, along with the offset vector given in (2.58), we pick Adjacency and pinning gain matrices given in (5.56). And, to test the stable independent of delay structure of the algorithms presented in Section 2.3, we added two seconds delay as a communication delay. By looking at Fig. 8, one can conclude that stable independent of delay property of the proposed algorithms, is verified.

### 2.5 Actual flight test results

This section reveals the flight test results obtained with single and multiple UAVs under the influence of time delays. We share the graphs of the desired and

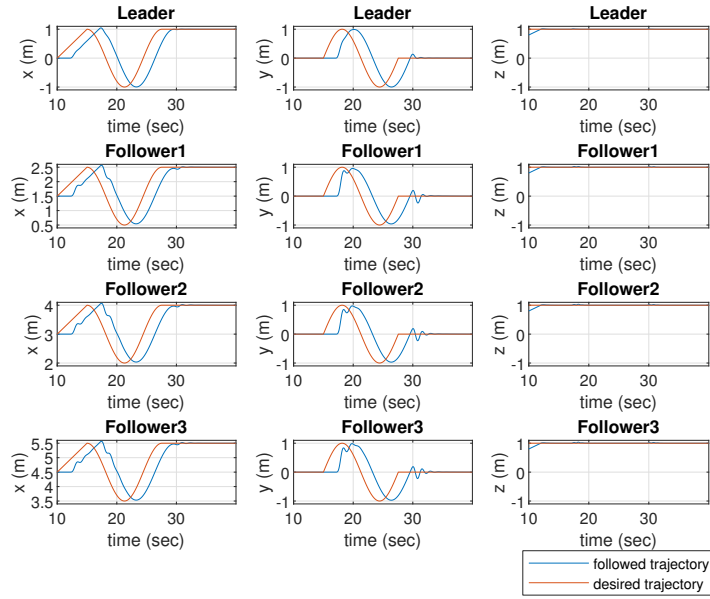


Figure 8: The path tracked by the UAVs when there exists two seconds delay and directed graph topology is used

followed trajectories when both the proposed methods in Section 2.2.3 and 2.3 is used.

### 2.5.1 Controller behavior with a single quadrotor

In this section, we present the performance of the backstepping control algorithm proposed in Section 3 by using both circular and figure-eight trajectories.

For the backstepping controller designed in Section 2.2.3, Fig. 9 and Fig. 10 shows the desired path of the UAV and the path traced by the UAV when both circular and eight-figure are desired trajectories. Notice that the tracking error is maintained inside the acceptable bounds, showing the performance of designed backstepping controller in terms of path following. Note that diagonal elements of  $\mathbf{K}_1$  and  $\mathbf{K}_2$ , are assigned respectively as 2, 2, 3 and 1.5, 1.5, 3.

To observe the network communication delay, we plot time vs desired and tracked positions in 3D space as shown in Fig. 11. Notice that communication delay

is about 2 seconds and tends to be commensurate through path following experiment. Moreover, when the quadrotor is following a nonlinear trajectory such as eight-figure and circular path, there exists time-delay between the desired position and followed position as shown in Fig. 11. This time delay is the summation of reaction time of UAV and the network delay caused by the local positioning system. That's why, the error seems to be bounded. However, after finishing the complete eight or circular figure trajectory, the error goes to zero since the desired position is constant at that time and quadrotor's position is the same as the desired position. This can be seen clearly from Fig. 9 and Fig. 10.

### 2.5.2 Controller behavior with multiple quadrotors

This section shows the formation control performance of the distributed backstepping control method given in Section 2.3. In the experiments of this section, the task of followers is to track the formation leader with a certain position offset. We show the leader and followers positions in the Fig. 12 and Fig. 13 while the leader

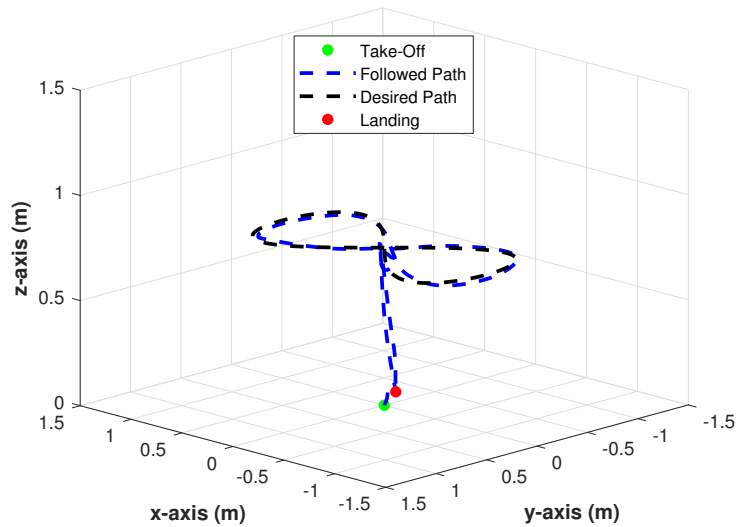


Figure 9: The path tracked by the UAV with the backstepping controller using 8-figure trajectory

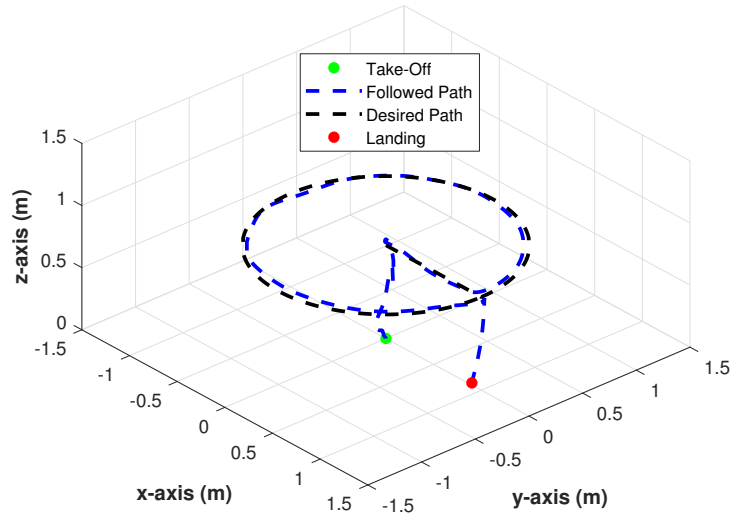


Figure 10: The path tracked by the UAV with the backstepping controller using a circular trajectory.

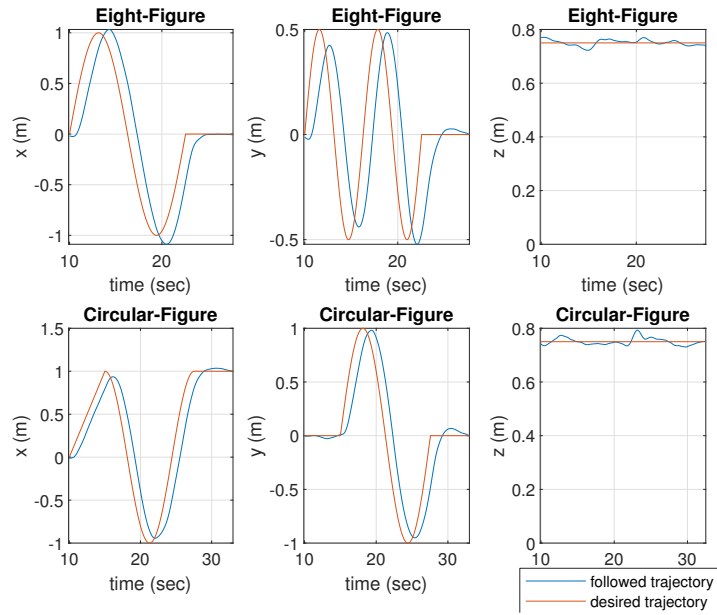


Figure 11: Controller behavior with time delay.

of the multi-agent system is following both eight-figure and circular trajectories re-

spectively. Adjacency and pinning gain matrices along with offset vector used in this actual hardware implementation are

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \mathbf{G} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \\ \underline{\Delta}^T &= \begin{bmatrix} 0 & -3 & 0 & 0 & -6 & 0 \end{bmatrix}. \end{aligned} \quad (2.60)$$

Fig. 14 and Fig. 15 show the desired path of the UAVs and the path tracked by the UAVs when the leader is following both circular and eight-figure type trajectories. Note that the delay experienced by each agent of the formation is slightly different than each other. However, if the position offset of the quadrotors gets bigger, agents of the formation would experience more distributed delays.

We record a video of the experiments described in this section, the interested reader can use the link <https://www.youtube.com/watch?v>

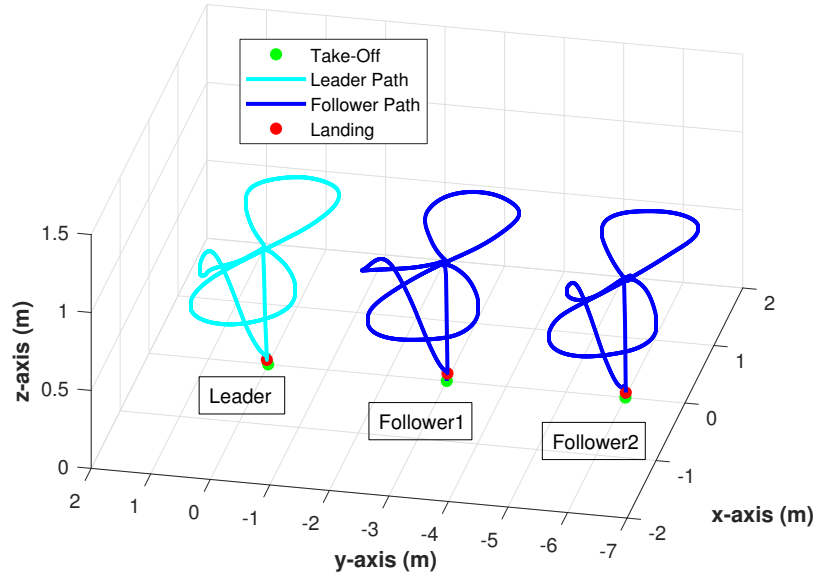


Figure 12: The path tracked by the UAVs with the distributed backstepping controller when the leader follows an eight-figure trajectory.

=rmY1LK42oPk' to have a visual understanding of the paper. Notice in the movie, formation control using the distributed backstepping method is influenced by the very strong wind effect that is produced by quadrotors themselves. This demonstrates the robustness of the proposed control method.

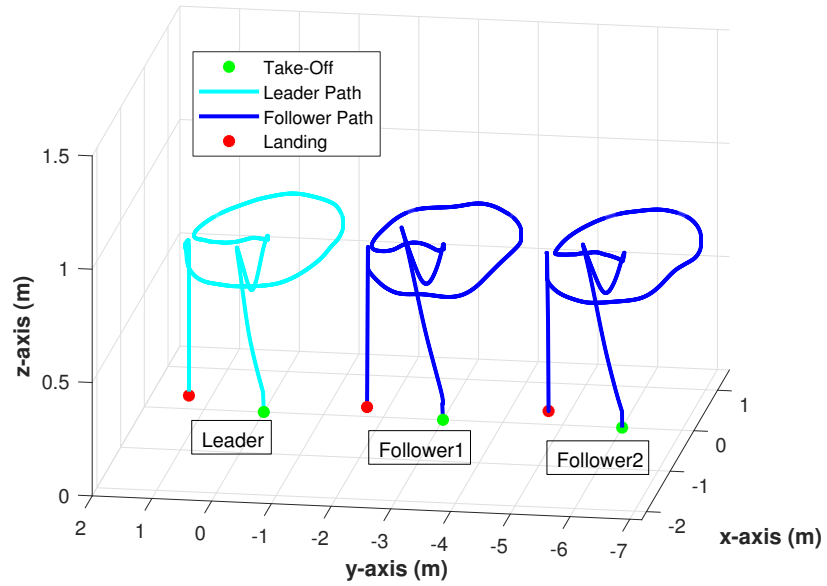


Figure 13: The path tracked by the UAVs with the distributed backstepping controller when the leader follows a circular trajectory.

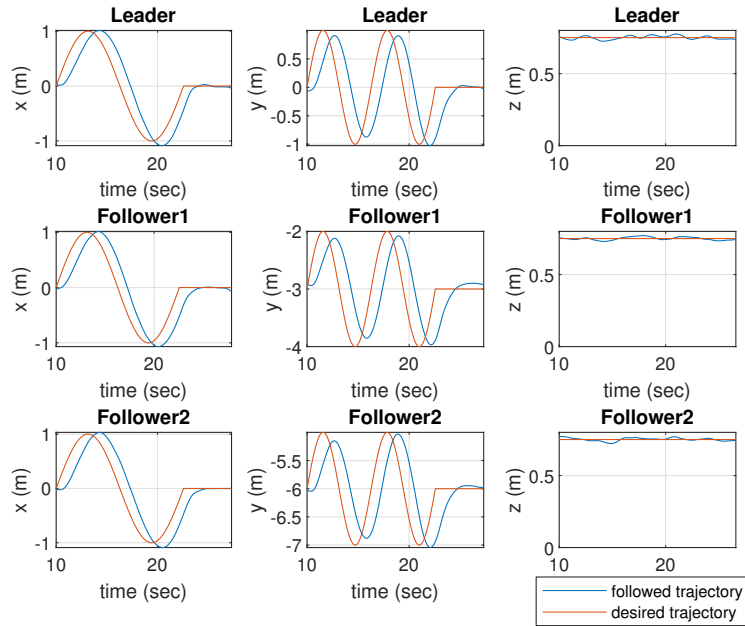


Figure 14: Observation of time delay graph when the leader follows an eight-figure trajectory.

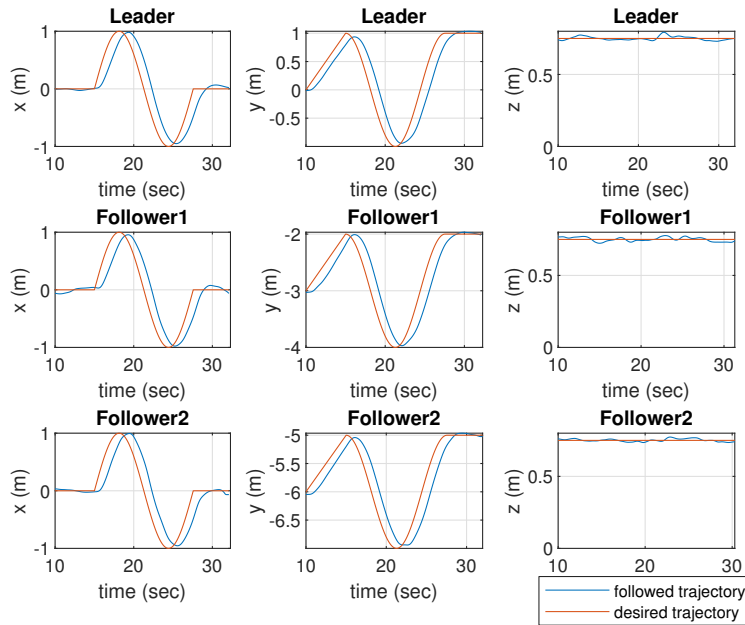


Figure 15: Observation of time delay graph when the leader follows a circular trajectory.



Optimal game theoretic solution of the pursuit-evasion intercept problem using  
on-policy reinforcement learning [54]

Y. Kartal, K.Subbarao, A. Dogan and F. Lewis, *Published in International Journal  
of Robust and Nonlinear Control*. Volume 31, Issue 16. 2021.

<https://doi.org/10.1002/rnc.5719>

## CHAPTER 3

### Optimal Game Theoretic Solution of the Pursuit-Evasion Intercept Problem Using On-Policy Reinforcement Learning

Inspired by the animal behaviors in hunting scenarios, the pursuit-evasion (PE) games have drawn great attention due to their applicability in areas such as missile guidance [55], collision avoidance systems [56] and controller designs [57]. The game of this kind is defined as a sub-category of differential game theory and provides the correct framework for the analysis of intercept problem and the choice of optimal policies for the agents involved in the two-player zero-sum (ZS) game.

Isaacs [7], founder of the differential game theory, initiated the development of strategic policies for both pursuer and evader in a PE problem. In Isaacs [58], the homicidal chauffeur game was analyzed in detail regarding players' speed and maneuverability capabilities. Bryson [59] introduced optimal feedback laws and demonstrated intercept strategies for players, by using the fixed final-time value function. Lewis et al. [60] made an extension of the Bellman equation, known as the Hamilton-Jacobi-Isaacs (HJI) equations to design  $H_\infty$  control, by employing the ZS games solutions. Moreover, works [12], [61], and [13] deal with linear-quadratic ZS games, in which their objective is to minimize the maximum norm of inputs and states, where the maximum is taken over the unknowns, such as disturbances. Hayoun et al. [62] reveal a set-based computing method for solving a general class of ZS Stackelberg differential games where the authors come up with a novel class of differential inequalities to get convex outer approximations of backward and forward reachable sets. Bhattacharya et al. [63] worked on a visibility-based PE game when the envi-

ronment contains a circular obstacle. Furthermore, Li et al. [64] and Liu et al. [65] developed an reinforcement learning (RL) algorithm to learn the Nash equilibrium solution for designing model-free controller by solving the game algebraic Riccati equation forward in time.

Applications of PE games involve the proximate satellite interception guidance strategies. The work [66] studied the intercept problem of satellites where both of the interceptor and target satellite can perform orbital maneuvers with limited thrusts. In the work [67], the authors analyzed same problem by establishing a local moving coordinate frame and simplifying dynamics of each player to the linear Clohessy-Wiltshire equations. Using the terminal time as a cost function, Gong et al. [68] derived sufficient conditions for capture in the PE problem based on the players' hyper-reachable domain. Note that the common point of these papers is to utilize prescribed terminal time on the construction of the game-theoretical cost function. Jagat et al. [69] proposes quadratic infinite-horizon cost functional for both players but the finite-time capture is not proven mathematically. Instead, simulations are provided to show that capture occurs in a finite-time. Carr et al. [70] employ semi-direct collocation nonlinear programming method to solve optimal actions for agents of the pursuit–evasion game. Authors solve the minimax problem by considering co-state dynamics and boundary conditions simultaneously for the dynamical models.

Standard solution to constrained PE game is to impose external velocity or acceleration constraints. Unfortunately, this leads to discontinuous saturated solutions that are difficult to analyze [59],[60].

Recent works by Hayoun et al. [62], Shaferman et al. [71] and Weintraub et al. [72] focus on the missile-target engagement where the PE problem is formulated as a differential game with an objective of optimizing the linear quadratic cost functional. The work [62] propose bounded maneuverability of the evader to prove the capture

in ZS game whereas Weintraub et al. [72] consider an engagement scenario by introducing the defense of a non-maneuverable agent. Further, this work [72] reveals the inclusion of altitude and dynamics in 3-dimensions, which is more realistic for the modeling of aerial engagements.

We sum up the contributions of this paper into four categories as:

- Firstly, a backstepping based velocity tracker is developed for PE games where the pursuer and evader both have arbitrary nonlinear dynamics. Taking a priori velocity constraints into account, a novel non-quadratic scalar functional is solved to obtain the smooth optimal velocity policies for each player in contrast to the standard discontinuous solutions.
- Secondly, with the detailed Lyapunov analysis, sufficient conditions are given for the case where capture must be attained in finite-time.
- The on-policy Integral Reinforcement Learning method is employed to solve the corresponding HJI equation and achieve the game optimal velocity policies for both pursuer and evader.
- Finally, the full rotational dynamics are added to extend the results to full nonlinear dynamical PE systems.

Rest of the paper is organized as follows. Section 3.2 reviews the exponentially stabilizing nonlinear backstepping control method to track given velocity trajectories for a generalized Newtonian system dynamics. Section 3.3 obtains optimal actions for the players by making use of the Pontryagin’s minimum principle and brings analysis of a Nash equilibrium in the PE game. Furthermore, having revealed the sufficient conditions for the asymptotic capture, we prove that PE game ends in a finite-time based on the derived sufficient conditions. Section 3.4 proposes an on-policy reinforcement learning algorithm for the solution of HJI equation and derives the proof of convergence to the optimal policies. Section 3.5 closes the backstepping

control loop by treating forces and/or moments as finalized inputs to the system and representing the attitude with unit quaternions to overcome the singularity problem of the Euler angles. Finally, the proposed control policies are illustrated via simulation results in Section 3.6.

### 3.1 Problem formulation and model description

We study the pursuit-evasion (PE) game for general Newtonian dynamics. A novel approach is given whereby we first design backstepping based velocity controllers for the pursuer and evader that guarantee a Nash solution to the PE game. We use a novel value function that ensures a solution under bounded velocities of the pursuer and the evader. This provides smooth solutions to the bounded velocity PE game in contrast to standard discontinuous solutions [59]. We conduct a Nash equilibrium analysis for a game of this kind. Further, we seek to obtain sufficient conditions for global exponential stability of the origin (equilibrium) using a rigorous Lyapunov analysis. Finally, we seek to derive conditions for a final time capture, and provide an upper bound on the time of capture.

The generalized translational and rotational dynamics for the pursuer and the evader can be modeled in their respective body frames of reference as

$$m^i \dot{\mathbf{v}}_B^i = m^i \mathbf{S}(\mathbf{w}_B^i) \mathbf{v}_B^i + \mathbf{N}^i \mathbf{f}_g^i + \mathbf{f}_B^i, \quad (3.1)$$

$$\mathbf{I}_B^i \dot{\mathbf{w}}_B^i = \mathbf{S}(\mathbf{w}_B^i) \mathbf{I}_B^i \mathbf{w}_B^i + \boldsymbol{\tau}_B^i \quad (3.2)$$

where the superscript  $i \in \{\mathbf{p}, \mathbf{e}\}$  with  $\mathbf{p}$  denoting the pursuer and  $\mathbf{e}$  denoting the evader respectively. Here,  $\mathbf{v}_B^i \in \mathbb{R}^3$ ,  $\mathbf{w}_B^i \in \mathbb{R}^3$  are the translational and angular velocities respectively, and  $\mathbf{f}_B^i \in \mathbb{R}^3$ ,  $\boldsymbol{\tau}_B^i \in \mathbb{R}^3$  are the control forces and moments respectively, in the body fixed reference frame. Further,  $\mathbf{I}_B^i \in \mathbb{R}^{3 \times 3}$  is the constant

nonsingular inertia matrix defined in the body frame and  $m^i$  is a scalar quantity that denotes the mass of players' rigid bodies. In addition,  $\mathbf{f}_g^i = [0 \ 0 \ m^i g]^T$  is the gravitational force vector whose components are written in the Inertial frame.  $\mathbf{S}(\mathbf{w}_B^i) \in \mathbb{R}^{3 \times 3}$  represents a skew-symmetric matrix form of the vector  $\mathbf{w}_B^i$ . Moreover,  $m^i \mathbf{S}(\mathbf{w}_B^i) \mathbf{v}_B^i$  and  $\mathbf{S}(\mathbf{w}_B^i) \mathbf{I}_B^i \mathbf{w}_B^i$  are due to the derivative of the body referenced linear and angular momentum of the vehicles relative to the Inertial frame.  $\mathbf{N}^i \in \mathbb{R}^{3 \times 3}$  is the rotation matrix from Inertial to body frame. Later in Section 3.5, we will call this Inertial frame as earth frame and give detailed explanation for the rotation matrix.

### 3.2 Developing velocity tracker using backstepping control method

In this section, we present an exponentially stabilizing backstepping control method to track given velocity trajectories. This velocity tracker is developed in this section, which uses only the translational dynamics (3.1). In Section 3.3, the velocity tracker is extended for PE games based on the translational dynamics (3.1) for both pursuer and evader. Then, in Section 3.5 we also consider rotational dynamics (3.2) to obtain general controllers for both velocity and attitude for pursuer and evader.

Note, we first derive the required velocity tracking control laws in the Inertial frame, and then subsequently Section 3.5 shows how they are realized using the dynamics in (3.1) and (3.2).

Thus, in the Inertial frame the translational dynamics is represented as,

$$m^i \dot{\mathbf{v}}^i = \mathbf{f}^i + \mathbf{f}_g^i \quad (3.3)$$

where  $\mathbf{v}^i \in \mathbb{R}^3$  is the velocity vector and  $\mathbf{f}^i = \mathbf{N}^{iT} \mathbf{f}_B^i$  is the control force, in the Inertial frame  $i \in \{\mathbf{p}, \mathbf{e}\}$ . Introducing a desired virtual force  $\mathbf{f}_d^i$ , to the system dynamics (3.3) we obtain

$$m^i \dot{\mathbf{v}}^i = \mathbf{f}_d^i + \mathbf{f}_g^i + \tilde{\mathbf{f}}^i \quad (3.4)$$

where  $\tilde{\mathbf{f}}^i = \mathbf{f}^i - \mathbf{f}_d^i$  is the difference of control and desired forces of the Newtonian system in 3-D.

Define velocity error as

$$\delta_v^i = \mathbf{v}_d^i - \mathbf{v}^i \quad (3.5)$$

where  $\mathbf{v}_d^i \in \mathbb{R}^3$  is the desired velocity designed for pursuer  $\mathbf{v}_d^p$  and evader  $\mathbf{v}_d^e$  in the next section. Take the derivative of (3.5) and substitute in (3.4) to obtain closed-loop velocity error dynamics as

$$m^i \dot{\delta}_v^i = m^i \dot{\mathbf{v}}_d^i - \mathbf{f}_g^i - \tilde{\mathbf{f}}^i - \mathbf{f}_d^i. \quad (3.6)$$

Then select ideal desired force as

$$\mathbf{f}_d^i = m^i \dot{\mathbf{v}}_d^i - \mathbf{f}_g^i + m^i \mathbf{K}^i \delta_v^i \quad (3.7)$$

where  $\mathbf{K}^i \in \mathbb{R}^{n \times n}$  is a positive-definite matrix. Substituting (3.7) in (3.6) yields

$$\dot{\delta}_v^i = -\mathbf{K}^i \delta_v^i - \frac{\tilde{\mathbf{f}}^i}{m^i}. \quad (3.8)$$

This enables us to derive exponential stability of the origin, as long as an admissible  $\mathbf{f}_d^i$  exists. In Section 3.5 we consider the rotational dynamics and show how to design the control force,  $\mathbf{f}^i$  and hence  $\mathbf{f}_B^i$  in (3.1) and (3.3) respectively, to make  $\tilde{\mathbf{f}}^i \rightarrow \mathbf{0}$  [27]. Then (3.8) shows that  $\delta_v^i \rightarrow \mathbf{0}$  exponentially.

**Remark 1.** *Tracking the vector quantity  $\mathbf{f}_d^i$  in (3.7) not only guarantees exponential stability of the equilibrium of (3.6) but also gives the desired attitude of the Newtonian system (3.3) so that it is aligned with the direction of  $\mathbf{f}_d^i$ .*

The next section deals with the derivation of optimal velocity trajectories  $\mathbf{v}_d^i$  for pursuer and evader, employed in (3.5). The design of the desired ideal forces  $\mathbf{f}_d^p$ ,  $\mathbf{f}_d^e$  is treated in Section 3.5.

### 3.3 Optimal game theoretic velocity generation for pursuit-evasion game

In this main section, we first propose a formulation of PE game and derive the optimal bounded desired velocity trajectories  $\mathbf{v}_d^p, \mathbf{v}_d^e$  in (3.5) for the players. Secondly, we conduct a Nash equilibrium analysis for the game and derive sufficient conditions for global exponential stability of the origin by rigorous Lyapunov analysis. Finally, conditions for finite-time capture and its upper bound are given.

#### 3.3.1 Pursuit-evasion game formulation

Assuming the players are governed by the velocity error dynamics (3.8), this section presents various definitions to develop the game-theoretically optimal solution of the PE game satisfying *velocity constraints* on the players. To simplify the notation, define desired velocity in (3.5) for the pursuer  $\mathbf{v}^p = \mathbf{v}_d^p$  and the evader  $\mathbf{v}^e = \mathbf{v}_d^e$ .

The following kinematic expressions enable us to derive desired velocities and thereby the forces (3.7) for pursuer and evader

$$\begin{aligned}\dot{\boldsymbol{\xi}}^p &= \mathbf{v}^p \\ \dot{\boldsymbol{\xi}}^e &= \mathbf{v}^e\end{aligned}\tag{3.9}$$



where  $\boldsymbol{\xi}^p \in \mathbb{R}^3$  and  $\boldsymbol{\xi}^e \in \mathbb{R}^3$  denote the 3-dimensional position vectors  $(x, y, z)$  of pursuer and evader respectively, which are defined with respect to Inertial frame. Hence  $\boldsymbol{v}^p \in \mathbb{R}^3$  and  $\boldsymbol{v}^e \in \mathbb{R}^3$  are desired velocity vectors of the pursuer and evader respectively. Note that (3.3) employs the translational velocity in the PE game. This allows analysis of ZS game for general nonlinear systems in Section 3.5.

Now, consider the following formulation for the zero-sum (ZS) PE game. Let the evader have an objective of maximizing the relative distance  $\boldsymbol{\delta} \in \mathbb{R}^3$ , defined as

$$\boldsymbol{\delta} = \boldsymbol{\xi}^p - \boldsymbol{\xi}^e, \quad (3.10)$$

whereas the pursuer tries to minimize (3.10). Moreover, let the velocities of both pursuer and evader be bounded by scalars  $|v_j^p| \leq \lambda^p; |v_j^e| \leq \lambda^e \forall j = 1, \dots, n$ . To satisfy these constraints, the value functional is defined as

$$V^{\pi^p, \pi^e}(\boldsymbol{\delta}) = \int_t^\infty \{\boldsymbol{\delta}^T \boldsymbol{Q} \boldsymbol{\delta} + U(\pi^p(\boldsymbol{\delta})) - U(\pi^e(\boldsymbol{\delta}))\} d\tau \quad (3.11)$$

where  $\boldsymbol{Q} \in \mathbb{R}^{n \times n}$  is a positive-definite matrix,  $\pi^p(\cdot)$  and  $\pi^e(\cdot)$  stand for the policies of pursuer and evader respectively in ZS game such that

$$\begin{aligned} \pi^p(\boldsymbol{\delta}) &\triangleq \boldsymbol{v}^p \\ \pi^e(\boldsymbol{\delta}) &\triangleq \boldsymbol{v}^e. \end{aligned} \quad (3.12)$$

Moreover,  $U(\boldsymbol{v}^i)$  (for  $i$  is either  $p$  or  $e$ ) is a generalized non-quadratic scalar functional [73], which ensures bounded velocities given by

$$U(\boldsymbol{v}^i) = 2 \int_{\mathbf{0}}^{\boldsymbol{v}^i} (\boldsymbol{\alpha}^{-1}(\boldsymbol{u}^i / \lambda^i))^T \boldsymbol{R}^i d\boldsymbol{u}^i, \quad (3.13)$$

$$\alpha^{-1}(\mathbf{u}^i/\lambda^i) = \left[ \alpha^{-1}(u_1^i/\lambda^i) \quad \dots \quad \alpha^{-1}(u_n^i/\lambda^i) \right]^T,$$

$$\mathbf{u}^i = \left[ u_1^i \quad \dots \quad u_n^i \right]^T, \mathbf{v}^i = \left[ v_1^i \quad \dots \quad v_n^i \right]^T$$

where  $\mathbf{R}^i \in \mathbb{R}^{n \times n}$  is a symmetric positive-definite matrix and  $\alpha(\cdot)$  is a bounded one-to-one smooth function i.e it belongs to  $C^\ell, \ell \geq 1$ . This is a monotonic odd function with its first derivative bounded by a constant. An example of  $\alpha(\cdot)$  is  $\tanh(\cdot)$  and throughout this paper, we use  $\tanh(\cdot)$ , which constrains the velocity to remain within predefined limits i.e  $|v_j^i| \leq \lambda, \forall j = 1, \dots, n$  and  $\forall i = p, e$ . In ZS PE games,  $\mathbf{R}^i$  plays a key role by restricting the rate of change of optimal velocities and hence constrains the accelerations of the each player.

The differential equivalent of (3.11) is the ZS game Bellman equation. Using (3.9), (3.10) and Leibniz's formula, the Bellman equation is obtained as

$$\begin{aligned} H(\boldsymbol{\delta}, \nabla V, \mathbf{v}^p, \mathbf{v}^e) &\equiv \boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\mathbf{v}^p) - U(\mathbf{v}^e) + \nabla V^T \dot{\boldsymbol{\delta}} \\ &\equiv \boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\mathbf{v}^p) - U(\mathbf{v}^e) + \nabla V^T (\mathbf{v}^p - \mathbf{v}^e) = 0 \end{aligned} \quad (3.14)$$

where  $\nabla V = \partial V^{\pi^p, \pi^e} / \partial \boldsymbol{\delta} \in \mathbb{R}^n$  is the gradient of value function (3.11), and  $H(\cdot)$  is the Hamiltonian.

To find the optimal policies  $\pi_*^i(\boldsymbol{\delta}) = \mathbf{v}_*^i$  (for  $i = p, e$ ) of players in the game, check stationarity conditions  $\partial H / \partial \mathbf{v}^p = \mathbf{0}$  and  $\partial H / \partial \mathbf{v}^e = \mathbf{0}$ . For the pursuer, applying Pontryagin's minimum principle to (3.14) yields

$$\frac{\partial H}{\partial \mathbf{v}^p} \equiv \frac{\partial U(\mathbf{v}^p)}{\partial \mathbf{v}^p} + \frac{\partial}{\partial \mathbf{v}^p} \{ \nabla V^T (\mathbf{v}^p - \mathbf{v}^e) \}. \quad (3.15)$$

Evaluating the derivatives at the right-hand-side of (3.15) using Leibniz's formula, and checking the stationarity condition  $\partial H/\partial \mathbf{v}^p = \mathbf{0}$  yields

$$2 \left( \tanh^{-1} \left( \frac{\mathbf{v}^{p*}}{\lambda^p} \right) \right)^T \mathbf{R}^p = -\nabla V^{*T}. \quad (3.16)$$

Then, the optimal policy for the pursuer using the definition (3.12) is obtained as

$$\pi^{p*}(\boldsymbol{\delta}) \triangleq \mathbf{v}^{p*} = -\lambda^p \tanh \left( \frac{1}{2} (\mathbf{R}^p)^{-1} \nabla V^* \right). \quad (3.17)$$

This velocity control bounded as required.

Likewise, one can follow the same steps to derive bounded optimal velocity policy for the evader as

$$\pi^{e*}(\boldsymbol{\delta}) \triangleq \mathbf{v}^{e*} = -\lambda^e \tanh \left( \frac{1}{2} (\mathbf{R}^e)^{-1} \nabla V^* \right). \quad (3.18)$$

Let  $V^*$  be the optimal value of (3.11) with the policies given in (3.17) and (3.18).

Then Hamilton-Jacobi-Isaacs (HJI) equation is obtained as

$$H(\boldsymbol{\delta}, \nabla V^*, \mathbf{v}_*^p, \mathbf{v}_*^e) \equiv \boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\mathbf{v}_*^p) - U(\mathbf{v}_*^e) + \nabla V^{*T} (\mathbf{v}_*^p - \mathbf{v}_*^e) = 0. \quad (3.19)$$

Note that the positive and negative definiteness of Hessians,  $\partial^2 H/\partial \mathbf{v}^{p2} > 0$  and  $\partial^2 H/\partial \mathbf{v}^{e2} < 0$ , indeed show that pursuer's optimal policy aims to minimize the Hamiltonian (3.14) whereas evader's aims to maximize. Therefore,  $(\mathbf{v}_*^p, \mathbf{v}_*^e)$  is the game-theoretic saddle point. Furthermore, this is a Nash equilibrium since the game is of type ZS and (3.11) is separable [59]. Rigorous analysis of this is shown in Theorem 1.

### 3.3.2 Proof of Nash equilibrium

In this section, we derive the value of PE game at Nash equilibrium. The following lemmas and corollary are necessary steps to prove that the Nash equilibrium is reached with policies (3.17) and (3.18).

**Lemma 1.** *Let  $V^{\pi^p, \pi^e}(\boldsymbol{\delta})$  be the corresponding solution of the Hamiltonian (3.14).*

*Then following equality holds*

$$\begin{aligned} H(\boldsymbol{\delta}, \nabla V, \mathbf{v}^p, \mathbf{v}^e) &= H(\boldsymbol{\delta}, \nabla V, \mathbf{v}_*^p, \mathbf{v}_*^e) + \nabla V^T((\mathbf{v}^p - \mathbf{v}_*^p) \\ &\quad + (\mathbf{v}_*^e - \mathbf{v}^e)) + U(\mathbf{v}^p) - U(\mathbf{v}_*^p) + U(\mathbf{v}_*^e) - U(\mathbf{v}^e). \end{aligned} \quad (3.20)$$

*Proof.* Adding and subtracting the terms  $U(\mathbf{v}_*^p)$ ,  $U(\mathbf{v}_*^e)$ ,  $\nabla V^T \mathbf{v}_*^p$  and  $\nabla V^T \mathbf{v}_*^e$  to Hamiltonian (3.14) yields

$$\begin{aligned} H(\boldsymbol{\delta}, \nabla V, \mathbf{v}^p, \mathbf{v}^e) &= \boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + \nabla V^T(\mathbf{v}_*^p - \mathbf{v}_*^e)U(\mathbf{v}_*^p) \\ &\quad - U(\mathbf{v}_*^e) + \nabla V^T((\mathbf{v}^p - \mathbf{v}_*^p) + (\mathbf{v}_*^e - \mathbf{v}^e)) \\ &\quad + U(\mathbf{v}^p) - U(\mathbf{v}_*^p) + U(\mathbf{v}_*^e) - U(\mathbf{v}^e), \end{aligned} \quad (3.21)$$

which completes the proof. □

**Lemma 2.** *Let  $V^{\pi^p, \pi^e}(\boldsymbol{\delta})$  be the corresponding solution of the Hamiltonian (3.14) and define  $V(\boldsymbol{\delta}(t_0))$  as the initial value of the game. Then following equality holds*

$$V^{\pi^p, \pi^e}(\boldsymbol{\delta}(t_0)) = \int_{t_0}^{\infty} H(\boldsymbol{\delta}, \nabla V, \mathbf{v}^p, \mathbf{v}^e) d\tau + V(\boldsymbol{\delta}(t_0)). \quad (3.22)$$

*Proof.* Assume that capture occurs in the interval  $t \in [t_0, \infty]$ , which implies  $\lim_{t \rightarrow \infty} V^{\pi^p, \pi^e}(\boldsymbol{\delta}(t)) =$

0. Then adding zero to (3.11) yields

$$\begin{aligned}
V^{\pi^p, \pi^e}(\boldsymbol{\delta}(t_0)) &= \int_{t_0}^{\infty} \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\pi^p(\boldsymbol{\delta})) - U(\pi^e(\boldsymbol{\delta}))\} d\tau + \int_{t_0}^{\infty} \dot{V}^{\pi^p, \pi^e} d\tau + V(\boldsymbol{\delta}(t_0)) \\
&= \int_{t_0}^{\infty} \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\mathbf{v}^p) - U(\mathbf{v}^e)\} d\tau + \int_{t_0}^{\infty} \nabla V^T(\mathbf{v}^p - \mathbf{v}^e) d\tau + V(\boldsymbol{\delta}(t_0)) \\
&= \int_{t_0}^{\infty} H(\boldsymbol{\delta}, \nabla V, \mathbf{v}^p, \mathbf{v}^e) d\tau + V(\boldsymbol{\delta}(t_0)). \tag{3.23}
\end{aligned}$$

This completes the proof.  $\square$

The next corollary extends the fact given in Lemma 1.

**Corollary 1.** *Suppose  $V^*$  satisfies the HJI equation (3.19). Then  $H(\boldsymbol{\delta}, \nabla V^*, \mathbf{v}_*^p, \mathbf{v}_*^e) = 0$  and (3.20) becomes*

$$H(\boldsymbol{\delta}, \nabla V^*, \mathbf{v}^p, \mathbf{v}^e) = \nabla V^{*T}((\mathbf{v}^p - \mathbf{v}_*^p) + (\mathbf{v}_*^e - \mathbf{v}^e)) + U(\mathbf{v}^p) - U(\mathbf{v}_*^p) + U(\mathbf{v}_*^e) - U(\mathbf{v}^e). \tag{3.24}$$

The next theorem derives the optimal value of the ZS game and proves Nash equilibrium reached.

**Theorem 1.** *Consider kinematic expressions for the players (3.9) with the value function given in (3.11). Let  $V^*$  be a positive definite smooth solution of HJI equation (3.19). Then,  $(\mathbf{v}_*^p, \mathbf{v}_*^e)$  given by (3.17), (3.18) is the Nash equilibrium and  $V^*(\boldsymbol{\delta}(t_0))$  is the value of PE game.*

*Proof.* Using the facts given in Lemma 2 and Corollary 1, (3.23) becomes

$$\begin{aligned}
V^{\pi^p, \pi^e}(\boldsymbol{\delta}(t)) &= \int_t^{\infty} \{\nabla V^{*T}((\mathbf{v}^p - \mathbf{v}_*^p) + (\mathbf{v}_*^e - \mathbf{v}^e)) \\
&\quad + U(\mathbf{v}^p) - U(\mathbf{v}_*^p) + U(\mathbf{v}_*^e) - U(\mathbf{v}^e)\} d\tau + V^*(\boldsymbol{\delta}(t_0)). \tag{3.25}
\end{aligned}$$

To prove  $(\mathbf{v}_*^p, \mathbf{v}_*^e)$  is the Nash equilibrium of the game, we need to show that when the pursuer adopts policy given in (3.17), the best action for the evader to maximize the value function (3.11) is  $\mathbf{v}^{e*}$ . Likewise, when the evader adopts policy given in (3.18), the best action for the pursuer to minimize the value function (3.11) is  $\mathbf{v}^{p*}$  i.e.

$$V^{\pi^{p*}, \pi^e}(\boldsymbol{\delta}(t)) \leq V^{\pi^{p*}, \pi^{e*}}(\boldsymbol{\delta}(t)) \leq V^{\pi^p, \pi^{e*}}(\boldsymbol{\delta}(t)). \quad (3.26)$$

Note that  $V^{\pi^{p*}, \pi^{e*}}(\boldsymbol{\delta}(t)) = V^*(\boldsymbol{\delta}(0))$  and call the integral term in (3.25) as  $\beta(V^{\pi^p, \pi^e})$ . Now we need to show  $\beta(V^{\pi^{p*}, \pi^e}) \leq 0$  and  $\beta(V^{\pi^p, \pi^{e*}}) \geq 0$  so that (3.26) holds. Then using (3.13), (3.17),(3.18) and (3.25) we obtain

$$\begin{aligned} \beta(V^{\pi^{p*}, \pi^e}) &= \int_t^\infty \{\nabla V^{*T}(\mathbf{v}_*^e - \mathbf{v}^e) + U(\mathbf{v}_*^e) - U(\mathbf{v}^e)\} d\tau \\ &= \int_t^\infty \{-2(\tanh^{-1}(\mathbf{v}^{e*}/\lambda^e))^T \mathbf{R}^e(\mathbf{v}_*^e - \mathbf{v}^e) + 2 \int_{\mathbf{v}^e}^{\mathbf{v}_*^e} (\tanh^{-1}(\mathbf{u}/\lambda^e))^T \mathbf{R}^e d\mathbf{u}\} d\tau. \end{aligned} \quad (3.27)$$

Now define  $\phi^T(\cdot) = \tanh^{-1}(\cdot)$  and note that  $\phi^T(\cdot)$  is monotonically increasing function in the interval  $[-\lambda^e, \lambda^e]$ . To complete the proof, first assume that  $\mathbf{v}^{e*} \geq \mathbf{v}^e$  and apply integral mean value theorem on (3.27)

$$\begin{aligned} \beta(V^{\pi^{p*}, \pi^e}) &= \int_t^\infty \{-2\phi(\mathbf{v}^{e*}/\lambda^e) \mathbf{R}^e(\mathbf{v}_*^e - \mathbf{v}^e) + 2 \int_{\mathbf{v}^e}^{\mathbf{v}_*^e} \phi(\mathbf{u}/\lambda^e) \mathbf{R}^e d\mathbf{u}\} d\tau \\ &\leq \int_t^\infty \{-2\phi(\mathbf{v}^{e*}/\lambda^e) \mathbf{R}^e(\mathbf{v}_*^e - \mathbf{v}^e) + 2\phi(\mathbf{v}^{e*}/\lambda^e) \mathbf{R}^e(\mathbf{v}_*^e - \mathbf{v}^e)\} d\tau = 0. \end{aligned} \quad (3.28)$$

Then assume that  $\mathbf{v}^{e*} < \mathbf{v}^e$  and again apply integral mean value theorem on (3.27)

$$\beta(V^{\pi^{p*}, \pi^e}) = \int_t^\infty \{2\phi(\mathbf{v}^{e*}/\lambda^e) \mathbf{R}^e(\mathbf{v}^e - \mathbf{v}^{e*}) - 2 \int_{\mathbf{v}_*^e}^{\mathbf{v}^e} \phi(\mathbf{u}/\lambda^e) \mathbf{R}^e d\mathbf{u}\} d\tau$$

$$\leq \int_t^\infty \{2\phi(\mathbf{v}^{e*}/\lambda^e)\mathbf{R}^e(\mathbf{v}^e - \mathbf{v}_*^e) - 2\phi(\mathbf{v}^{e*}/\lambda^e)\mathbf{R}^e(\mathbf{v}^e - \mathbf{v}_*^e)\}d\tau = 0, \quad (3.29)$$

which shows that  $\beta(V^{\pi^{p*}, \pi^e}) \leq 0$ . The same procedure can be performed to show  $\beta(V^{\pi^p, \pi^{e*}}) \geq 0$ . Then the inequality given in (3.26) is verified, which implies that  $(\mathbf{v}_*^p, \mathbf{v}_*^e)$  is the Nash equilibrium and  $V^*(\boldsymbol{\delta}(t_0))$  is the value of the PE game.  $\square$

### 3.3.3 Stability and finite-time capture analysis

This section first reveals the sufficient conditions for the asymptotic capture of the evader by the pursuer. Then, by making use of these conditions, derives the globally exponential stability of the origin. Finally it is shown that under certain conditions, finite-time capture is ensured.

Before developing analysis for the asymptotic capture, let  $\mathbf{R}^i$  in (3.17) and (3.18) be a diagonal matrix with elements of  $r_j^i > 0, \forall j \in \{1, 2, 3\}$  and  $\forall i = p, e$ . This enables us to simplify the analysis that will be developed in the rest of the paper. Employing this assumption, the next theorem shows the sufficient conditions for the asymptotic capture in ZS PE games.

**Theorem 2.** *Consider kinematic expressions for the players (3.9) with the value function given in (3.11). Then the equilibrium of tracking error dynamics  $\dot{\boldsymbol{\delta}} = \mathbf{v}^{p*} - \mathbf{v}^{e*}$ , is asymptotically stable point with candidate Lyapunov function  $L(\boldsymbol{\delta}) = V^{\pi^{p*}, \pi^{e*}}(\boldsymbol{\delta})$ . The sufficient conditions for asymptotic capture are  $\lambda^p > \lambda^e$  and  $r_{e_i} \geq r_{p_i}, \forall i = 1, \dots, n$ .*

*Proof.* Since  $V^{\pi^{p*}, \pi^{e*}}(\boldsymbol{\delta})$  does not depend on the time explicitly, equality  $\dot{L}(\boldsymbol{\delta}) = \nabla L^T \dot{\boldsymbol{\delta}}$  holds. By (3.14), derivative of the Lyapunov function,  $\dot{L}(\boldsymbol{\delta})$  is obtained as

$$\dot{L}(\boldsymbol{\delta}) = -\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} - U(\mathbf{v}^{p*}) + U(\mathbf{v}^{e*}). \quad (3.30)$$

Assumption of the equality  $\mathbf{R}^p = \mathbf{R}^e$ , implies the pursuer and evader are moving in the same direction by (3.17) and (3.18). For intercept, the position of the pursuer and evader must be equal. To meet this criteria, we propose  $\lambda^p > \lambda^e$  so that asymptotic capture occurs as the row elements of optimal actions satisfy  $|v^{p*}|_i > |v^{e*}|_i \forall i = 1, \dots, n$ . Furthermore, taking (3.30) into account, the condition of  $\mathbf{R}^p = \mathbf{R}^e$  is relaxed as  $\mathbf{R}^e \geq \mathbf{R}^p$  since proposition  $\lambda^p > \lambda^e$  implies  $U(\mathbf{v}^{p*}) \geq U(\mathbf{v}^{e*})$ ,  $\dot{L}(\boldsymbol{\delta})$  becomes strictly negative definite. Then sufficient conditions for the asymptotic capture is proved to be  $\lambda^p > \lambda^e$  and  $r_{e_i} \geq r_{p_i} \forall i = 1, \dots, n$ .  $\square$

**Remark 2.** *Asymptotic capture in Theorem 2 can be strengthened to finite-time capture with the assumption that players involved in the game satisfy the sufficient conditions derived in the proof of Theorem 2. See Lemma 3.*

Following theorem extends the Theorem 2 to exponential stability of the origin.

**Theorem 3.** *Consider sufficient conditions and Lyapunov function,  $L(\boldsymbol{\delta})$  given in Theorem 2. Then, there exists positive scalars  $c_1, c_2$  and  $\epsilon$ , which satisfies*

$$\begin{aligned} c_1 \|\boldsymbol{\delta}\|_2^2 &\leq L(\boldsymbol{\delta}) \leq c_2 \|\boldsymbol{\delta}\|_2^2 \\ \dot{L}(\boldsymbol{\delta}) &\leq -\epsilon L(\boldsymbol{\delta}), \end{aligned} \tag{3.31}$$

*which implies that the origin is an exponentially stable equilibrium. Furthermore, radially unboundedness of the  $L(\boldsymbol{\delta})$  implies the globally exponentially stability of the origin [74], which is an essential result as the initial positional offset between the pursuer and evader should not be problem to prove the capture in PE game.*

*Proof.* The inequality  $U(\mathbf{v}^{p*}) \geq U(\mathbf{v}^{e*})$  by Theorem 2 and the strict convexity of  $U(\mathbf{v}_i)$  (for  $i = p, e$ ), imply the existence of positive scalars  $c_1$  and  $c_2$  [75]. Now, define



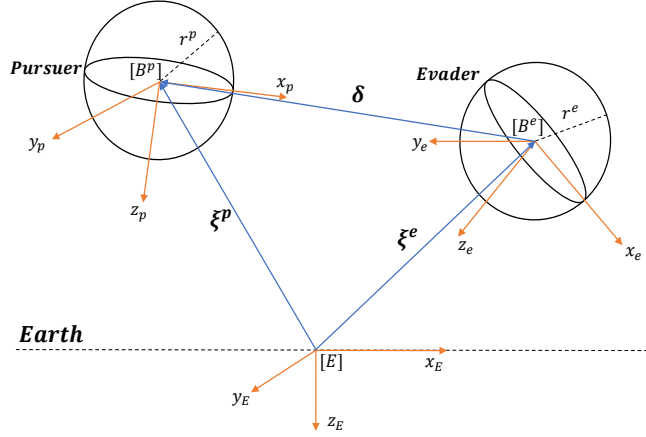


Figure 16: Sphere of collisions for players and their frames in 3-dimensions used for finite-time capture analysis.

convex function  $U_s(\boldsymbol{\delta})$  that satisfies the inequality  $U_s(\boldsymbol{\delta}) \leq U(\mathbf{v}^{p*}) - U(\mathbf{v}^{e*})$ . Using this and (3.30), the following inequality is derived as

$$\dot{L}(\boldsymbol{\delta}) \leq -\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} - U_s(\boldsymbol{\delta}). \quad (3.32)$$

Substituting (3.32) in (3.11) with the optimal policies (3.17) and (3.18), results in

$$L(\boldsymbol{\delta}) \leq \int_t^\infty \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U_s(\boldsymbol{\delta})\} d\tau, \quad (3.33)$$

which stands for the proof of  $\dot{L}(\boldsymbol{\delta}) \leq -\epsilon L(\boldsymbol{\delta})$  for sufficiently small  $\epsilon$ , which completes the proof.  $\square$

Notice that PE game given in Section 3.3 is formulated by treating the players as unit masses since the kinematic expressions (3.9) is employed in the value function (3.11). In Section 3.5, we consider full nonlinear dynamics (3.1), (3.2). Now, consider

the volume of pursuer and evader in 3-dimensional space and let the pursuer and evader have a sphere of collision with radius  $r^p$  and  $r^e$  respectively, as illustrated in Fig. 26. Then capture occurs when the distance between the center of masses of players is less than  $r^p + r^e$ . With this in mind, the next main lemma proves that the capture of evader by pursuer indeed occurs in finite-time in PE game.

**Lemma 3.** *There exists an upper-bound for the capture time in PE game when the conditions  $\lambda^p > \lambda^e$  and  $r_{e_i} \geq r_{p_i} \forall i = 1, \dots, n$  derived in Theorem 2 are satisfied. This also implies that the PE game ends in finite-time as required.*

*Proof.* The globally exponentially stability of the origin derived in Theorem 3 implies the equation of positional offset between the players is in the form of

$$\|\delta(t)\|_2 \leq b_1 \|\delta(t_0)\|_2 e^{-b_2(t-t_0)} \quad \forall t > t_0 \quad (3.34)$$

where  $b_i$  is a positive scalar  $\forall i = 1, 2$ . Then the upper bound for capture time  $t_c$  is derived as

$$t_c \leq t_0 + \frac{1}{b_2} \log \left( \frac{b_1 \|\delta(t_0)\|_2}{r^e + r^p} \right) \quad (3.35)$$

where  $\log(\cdot)$  is a natural logarithm function and this completes the proof.  $\square$

**Remark 3.** *It is seen that for finite-time capture, the velocity bound  $\lambda^p$  for the pursuer must be greater than the velocity bound  $\lambda^e$  on the evader. Moreover, the sufficient condition on weights (3.13) is found as  $r_{e_i} \geq r_{p_i} \forall i = 1, \dots, n$ . Note that capture time is also studied for multi-agent systems in the work [76] by assuming the players are using their maximum efforts. In Lemma 3, we showed that capture time is upper bounded under certain conditions even the players are not using their maximum efforts.*

### 3.4 Online solution of HJI equation using integral reinforcement learning (IRL)

The PE game formulation in Section 3.3 requires the generation of velocity set-points online and in real-time for both agents of the game. With this in mind, we employ the following synchronous IRL algorithm [8] to solve the HJI equation (3.19) in real-time and hence, reach the Nash equilibrium velocity policies  $(\mathbf{v}_*^p, \mathbf{v}_*^e)$  online by observing measured data. In the work [28], it is emphasized that persistence of excitation condition must be satisfied so that the IRL algorithm converges. This is achieved in most applications [8]-[9] by adding small probing noise. In our case, the persistence of excitation for IRL to work is satisfied inherently until capture occurs, at which time the game ends.

The tracking HJI equation (3.19) is nonlinear in the value function gradient  $\nabla V^*$ , and non-quadratic partial differential equation that is extremely difficult to solve. However, (3.36) can be solved for the value function and its gradient by collecting position data over some interval  $[t, t + T]$ . Therefore, finding the value of game optimal velocity policies by solving (3.36) is easier than solving (3.19). This is the motivation of introducing an iterative algorithm for approximating the tracking HJI solution, which is necessary to evaluate game optimal velocity policy for pursuer (3.17), and evader (3.18).

#### 3.4.1 Policy iteration solution for PE game

In this section, we present a policy iteration algorithm that avoids solution of (3.19) and also does not require knowledge of the system dynamics. The following lemma enables us to recognize IRL form of the value function (3.11).

**Lemma 4.** Let  $V^{\pi^p, \pi^e}(\boldsymbol{\delta})$  be the corresponding solution of the Bellman equation (3.14). Then, the value function (3.11), can be written in the IRL form as

$$V^{\pi^p, \pi^e}(\boldsymbol{\delta}(t)) = \int_t^{t+T} \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\pi^p(\boldsymbol{\delta})) - U(\pi^e(\boldsymbol{\delta}))\} d\tau + V^{\pi^p, \pi^e}(\boldsymbol{\delta}(t+T)). \quad (3.36)$$

*Proof.* The equality  $\dot{V}^{\pi^p, \pi^e} = -\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} - U(\pi^p(\boldsymbol{\delta})) + U(\pi^e(\boldsymbol{\delta}))$  holds by the differentiation of ZS game Bellman equation (3.14). Then integrating both sides from  $t$  to  $t+T$ , results in

$$\int_t^{t+T} \dot{V} d\tau = - \int_t^{t+T} \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\pi^p(\boldsymbol{\delta})) - U(\pi^e(\boldsymbol{\delta}))\} d\tau, \quad (3.37)$$

which verifies (3.36).  $\square$

The online policy-iteration Algorithm 1 performs a sequence of four-step iterations to find the optimal control policies for players. Notice that these policies stand for the optimal desired velocities, which are employed in (3.5). Furthermore, they are also Nash equilibrium velocity policies by Theorem 1.

1. Select any policy  $\pi_0^p$  and  $\pi_0^e$  for the players
2. Policy evaluation

$$V^{\pi_j^p, \pi_j^e}(\boldsymbol{\delta}(t)) = \int_t^{t+T} \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\pi_j^p(\boldsymbol{\delta})) - U(\pi_j^e(\boldsymbol{\delta}))\} d\tau + V^{\pi_j^p, \pi_j^e}(\boldsymbol{\delta}(t+T)). \quad (3.38)$$

3. Policy improvement

$$\begin{aligned} \pi_{j+1}^p(\boldsymbol{\delta}) &= -\lambda^p \tanh\left(\frac{1}{2}(\mathbf{R}^p)^{-1} \nabla V^{\pi_j^p, \pi_j^e}\right), \\ \pi_{j+1}^e(\boldsymbol{\delta}) &= -\lambda^e \tanh\left(\frac{1}{2}(\mathbf{R}^e)^{-1} \nabla V^{\pi_j^p, \pi_j^e}\right). \end{aligned} \quad (3.39)$$

4. On convergence stop; else go to step 2.  $\square$

**Algorithm 1:** Online policy-iteration algorithm

Notice that the position data of each player is collected through each iteration over the period  $T$ . The proof of convergence of Algorithm 1 to the optimal policies is shown in the following theorem.

**Theorem 4.** *Using the temporal difference (TD) learning method, Algorithm 1 converges to the Nash value  $V^*(\boldsymbol{\delta}(t_0))$  and Nash equilibrium policies  $(\pi^{p*}, \pi^{e*})$ , which optimizes velocity trajectories for the players in a game theoretic manner.*

*Proof.* First, evaluate the value function  $V^{\pi^{pj}, \pi^{ej}}(\boldsymbol{\delta}(t))$ , which solves the (3.38) by TD method. Then by Theorem 1, Isaacs' condition is derived as

$$H(\boldsymbol{\delta}, \nabla V, \mathbf{v}_*^p, \mathbf{v}_*^e) \leq H(\boldsymbol{\delta}, \nabla V, \mathbf{v}_*^p, \mathbf{v}_*^e) \leq H(\boldsymbol{\delta}, \nabla V, \mathbf{v}^p, \mathbf{v}^e). \quad (3.40)$$

Noting  $\beta(V^{\pi^{p*}, \pi^{e*}}) \leq 0$  and  $\beta(V^{\pi^p, \pi^e}) \geq 0$  is proved in the Theorem 1, the uniform convergence of Algorithm 1 immediately follows from Dini's theorem as reinforcement  $H(\boldsymbol{\delta}, \nabla V^j, \mathbf{v}^p, \mathbf{v}^e)$  converges to  $H(\boldsymbol{\delta}, \nabla V^*, \mathbf{v}_*^p, \mathbf{v}_*^e) = 0$  by Corollary 1. Moreover, due to the uniqueness of the value function (3.11), it follows that  $\lim_{j \rightarrow \infty} V^{\pi_j^p, \pi_j^e}(\boldsymbol{\delta}(t)) = V^*(\boldsymbol{\delta}(t_0))$ .  $\square$

### 3.4.2 Value function approximation to find game optimal pursuer and evader velocity policies

This section presents a critic neural network structure for policy-evaluation step in Algorithm 1.

**Remark 4.** *The IRL method given in Algorithm 1 requires the value function approximation (VFA), which can be achieved in a least-squares sense that is also known as single hidden layer critic Neural Network (NN). We employ this technique as in [8] that guarantees the successive least-squares iterations converge to the optimal value function of the HJI equation (3.19), and hence  $\nabla V^*$ .*

**Remark 5.** Note that the pair  $(\mathbf{v}_*^p, \mathbf{v}_*^e)$  stands for the Nash equilibrium by Theorem 1, thereby the Algorithm 1 converges to optimal actions for both players. Unlike the works [9] and [77] that use the IRL technique to reach min or max point of the value functional, we employ this technique to converge game theoretic saddle point by using the Isaacs' condition derived in Theorem 4. In addition, the system dynamics (3.1) does not appear in the value functional, which implies that we do not need to implement actor NN [77] [78] and the solution of HJI (3.19) can be obtained by using only critic NN, see [79].

By Remarks 4 and 10, we approximate the game optimal value functional in step 2 of Algorithm 1 using Weierstrass approximator such that

$$\begin{aligned}\hat{V}(\boldsymbol{\delta}) &= \hat{\mathbf{W}}^T \Phi(\boldsymbol{\delta}), \\ \nabla \hat{V} &= \Phi(\boldsymbol{\delta})^T \hat{\mathbf{W}}\end{aligned}\tag{3.41}$$

where  $\Phi(\boldsymbol{\delta}) \in \mathbb{R}^{nk}$  is the  $k$ -times concatenated basis function vector,  $n = 3$  as  $\boldsymbol{\delta} \in \mathbb{R}^3$ , and  $\hat{\mathbf{W}}$  is a critic NN weight vector to be determined. Using (3.41), the *policy evaluation* step of the IRL Algorithm 1 can be re-written as

$$e_b = \hat{\mathbf{W}}^T \Delta \Phi(\boldsymbol{\delta}) - \kappa(t)\tag{3.42}$$

where  $e_b$  is the continuous-time counterpart residual error of the TD,  $\Delta \Phi(\boldsymbol{\delta}) = \Phi(\boldsymbol{\delta}(t)) - \Phi(\boldsymbol{\delta}(t + T))$ , and reinforcement

$$\kappa(t) = \int_t^{t+T} \{\boldsymbol{\delta}^T \mathbf{Q} \boldsymbol{\delta} + U(\pi^p(\boldsymbol{\delta})) - U(\pi^e(\boldsymbol{\delta}))\} d\tau.\tag{3.43}$$

Therefore, (3.42) implies that the problem of solving the HJI equation is converted to tuning the critic NN weights such that  $e_b$  to be minimized. Now, to adjust these weights, the following objective function is employed

$$E_b = \frac{1}{2}e_b^2. \quad (3.44)$$

Then, the TD gradient descent algorithm [77] to minimize  $e_b$  is obtained by using the chain rule

$$\dot{\mathbf{W}} = -\frac{\alpha_L \Delta \Phi(\boldsymbol{\delta})}{(1 + \Delta \Phi(\boldsymbol{\delta})^T \Delta \Phi(\boldsymbol{\delta}))^2} e_b \quad (3.45)$$

where  $\alpha_L > 0$  is the learning rate. The proof of convergence of critic NN weights is shown in the Theorem 3 of Modares et al. [77].

### 3.5 Generalized rotational dynamics of the pursuer and evader

The analysis in the preceding sections has shown how to derive velocity tracker for the PE game given velocity dynamics (3.3). In this section, we analyze the general rotational dynamics (3.2) that are coupled to (3.1), and hence (3.3). We first derive the desired attitude of the system by using the  $Z$ - $Y$ - $X$  Euler angle rotation matrix from [E] (earth frame) to [B<sup>*i*</sup>] (body frames of pursuer or evader) as shown in Fig. 26, and desired force vector  $\mathbf{f}_d^i$  in (3.7). Then, by the analysis developed on the desired Euler angles, we propose the desired attitude representation with unit quaternions to overcome the singularity problem of the Euler angles. Lastly, by treating forces and/or moments as final inputs to the Newtonian system, we close the backstepping control loop to track desired force vector  $\mathbf{f}_d^i$  in (3.7). Note that in this section,  $i$  represents either  $p$  or  $e$ .

Assume that the gravity  $g$  is constant and the Earth is flat in the 3-dimensional space as illustrated in the Fig. 26. Then, the vehicle carrier frame is aligned with the body frame  $[B^i]$ . Thereby the rotation matrix from  $[E]$  to  $[B^i]$  frames shown in Fig. 26, can be given in terms of the Euler angles as

$$\mathbf{N}(\boldsymbol{\eta}^i) = \begin{bmatrix} c\theta^i c\psi^i & c\theta^i s\psi^i & -s\theta^i \\ -c\varphi^i s\psi^i + s\varphi^i s\theta^i c\psi^i & c\varphi^i c\psi^i + s\varphi^i s\theta^i s\psi^i & s\varphi^i c\theta^i \\ s\varphi^i s\psi^i + c\varphi^i s\theta^i c\psi^i & -s\varphi^i c\psi^i + c\varphi^i s\theta^i s\psi^i & c\varphi^i c\theta^i \end{bmatrix} \quad (3.46)$$

where  $c$  and  $s$  refers to cosine and sine respectively, and  $\boldsymbol{\eta}^i = [\psi^i \ \theta^i \ \varphi^i]^T$  is the Euler angle vector. Note that  $\mathbf{N}(\boldsymbol{\eta}^i)$  belongs to the special orthogonal group and is of rank 3, or  $SO(3)$ , whose determinant is equal to 1.

Assuming the direction of the thrust force to be along the nose of players' bodies or positive  $x_i$ -axis ( $\forall i = p, e$ ). This enables us to write that the desired force vector is indeed in the form of  $\mathbf{f}_{B_d}^i = [\mu_d^i \ 0 \ 0]^T$ , whose components written in  $[B^i]$ . Using (3.7) and expressing the desired force  $[\mu_d^i \ 0 \ 0]^T$  in  $[E]$  by  $\mathbf{f}_d^i = \mathbf{N}^T(\boldsymbol{\eta}_d^i) \mathbf{f}_{B_d}^i$ , following relation is derived

$$\mathbf{f}_d^i = \begin{bmatrix} f_{x_d}^i \\ f_{y_d}^i \\ f_{z_d}^i \end{bmatrix} = \mathbf{N}^T(\boldsymbol{\eta}_d^i) \begin{bmatrix} \mu_d^i \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \mu_d^i (c\theta_d^i c\psi_d^i) \\ \mu_d^i (c\theta_d^i s\psi_d^i) \\ \mu_d^i (-s\theta_d^i) \end{bmatrix} \quad (3.47)$$

where  $\boldsymbol{\eta}_d^i = [\psi_d^i \ \theta_d^i \ \varphi_d^i]^T$  is the desired Euler angle vector.

Then, (3.47) can be solved for desired attitude angles  $\theta_d$ ,  $\psi_d$  and  $\mu_d$  as

$$\theta_d^i = -\tan^{-1} \left( \frac{f_{z_d}^i}{f_{x_d}^i \cos\psi_d^i + f_{y_d}^i \sin\psi_d^i} \right), \quad (3.48)$$



$$\psi_d^i = \tan^{-1} \left( \frac{f_{y_d}^i}{f_{x_d}^i} \right), \quad (3.49)$$

$$\mu_d^i = \sqrt{f_{x_d}^i{}^2 + f_{y_d}^i{}^2 + f_{z_d}^i{}^2}. \quad (3.50)$$

Note that  $\varphi_d^i$  can be arbitrarily prescribed. However, (3.48)-(3.50) assumes that the equality conditions  $f_{x_d}^i = 0$ ,  $f_{y_d}^i = 0$  cannot occur simultaneously since (3.48) and (3.49) become indefinite. This singularity problem is also known as gimbal lock, which is associated with  $\theta_d^i = \pi/2$ .

To avoid gimbal lock, define the following unit quaternion representation

$$\mathbf{q}^i = [q_0^i \ q_1^i \ q_2^i \ q_3^i]^T = [q_0^i; \mathbf{q}_v^i] \quad (3.51)$$

$$q_0^i = \cos \phi^i / 2 \quad (3.52)$$

$$\mathbf{q}_v^i = \mathbf{k}^i \sin \phi^i / 2 \quad (3.53)$$

where  $\phi^i$  is the rotation about equivalent axis  $\mathbf{k}^i$ . Moreover, the kinematics equation for unit quaternion is

$$\dot{\mathbf{q}}^i = \frac{1}{2} \mathbf{J}^T(\mathbf{q}^i) \mathbf{w}_B^i \quad (3.54)$$

where  $\mathbf{J}(\mathbf{q}^i) \in \mathbb{R}^{3 \times 4}$  satisfies the equalities  $\mathbf{J}(\mathbf{q}^i) \mathbf{J}^T(\mathbf{q}^i) = \mathbf{I}_{3 \times 3}$ ,  $\mathbf{J}(\mathbf{q}^i) \mathbf{q}^i = \mathbf{0}$ , and can be expressed as

$$\mathbf{J}(\mathbf{q}^i) = [-\mathbf{q}_v^i \quad \mathbf{S}(\mathbf{q}_v^i) + q_0^i \mathbf{I}_{3 \times 3}] \quad (3.55)$$

where  $\mathbf{S}(\mathbf{q}_v^i) = \begin{bmatrix} 0 & q_3 & -q_2 \\ -q_3 & 0 & q_1 \\ q_2 & -q_1 & 0 \end{bmatrix}$ .

Then, the rotation matrix from  $[B^i]$  to  $[E]$  in terms of the unit quaternion (3.51) is given by

$$\mathbf{N}^T(\mathbf{q}^i) = \mathbf{I}_{3 \times 3} - 2q_0^i \mathbf{S}(\mathbf{q}_v^i) + 2\mathbf{S}^2(\mathbf{q}_v^i), \quad (3.56)$$

which is also known as *Rodrigues* formula. The following set of equations can be obtained by substituting the rotation matrix with the argument  $\mathbf{q}_d^i$  (3.56) into (3.47) along with selected  $\varphi_d^i$

$$\begin{bmatrix} f_{x_d}^i \\ f_{y_d}^i \\ f_{z_d}^i \end{bmatrix} = \mu_d^i \begin{bmatrix} 1 + 2(-q_{2_d}^i{}^2 - q_{3_d}^i{}^2) \\ 2q_{0_d}^i q_{3_d}^i + q_{1_d}^i q_{2_d}^i \\ -2q_{0_d}^i q_{2_d}^i + q_{1_d}^i q_{3_d}^i \end{bmatrix} \quad (3.57)$$

$$\varphi_d = \tan^{-1} \left( \frac{2(q_{0_d}^i q_{1_d}^i + q_{2_d}^i q_{3_d}^i)}{1 - 2(q_{1_d}^i{}^2 + q_{2_d}^i{}^2)} \right).$$

Notice that  $\mathbf{f}_d^i = [f_{x_d}^i \ f_{y_d}^i \ f_{z_d}^i]^T$ ,  $\varphi_d^i$  and  $\mu_d^i$  are known by (3.7) and (3.48)-(3.50). Thence, (3.57) can be solved for the desired unit quaternion  $\mathbf{q}_d^i = [q_{0_d}^i \ q_{1_d}^i \ q_{2_d}^i \ q_{3_d}^i]^T$  as (3.57) represents four equations with four unknowns, which are the elements of  $\mathbf{q}_d^i$ . Further substitute  $\mathbf{q}_d^i$  into kinematics equation (3.54) to find the desired angular velocity  $\mathbf{w}_{B_d}^i$  such that

$$\mathbf{w}_{B_d}^i = 2\mathbf{J}(\mathbf{q}_d^i)\dot{\mathbf{q}}_d^i. \quad (3.58)$$

**Remark 6.** For any Newtonian system (3.1) or (3.3) and (3.2), we know that forces and moments are coupled to each other, which implies that  $\boldsymbol{\tau}_B^i$  is required to be compatible with the selected desired force  $\mathbf{f}_d^i$  in (3.7).

Then applying the dynamic inversion technique,  $\boldsymbol{\tau}_B^i$  is given using (3.58) as

$$\boldsymbol{\tau}_B^i = \mathbf{I}_B^i \dot{\mathbf{w}}_{B_d}^i - \mathbf{S}(\mathbf{w}_{B_d}^i) \mathbf{I}_B^i \mathbf{w}_{B_d}^i. \quad (3.59)$$

Notice that we treat  $\boldsymbol{\tau}_B^i$  as final input for the general rotational dynamics (3.2). Consequently, we will not develop further analysis by giving location of thrusters and actuators, which is a control allocation problem and out of scope of this paper. Interested reader can check our work [27] to examine how to generate  $\boldsymbol{\tau}_B^i$  for the quadrotors.

### 3.6 Implementation on dynamic system

This section reveals the simulation results of ZS PE game with different scenarios. First, we consider when both the pursuer and evader follows their game optimal velocities given in (3.17) and (3.18) respectively. Then, we show the scenario in which the pursuer tracks its game optimal velocity (3.17) whereas the evader adopts a sub-optimal velocity policy.

In order to model the constrained optimal velocity trajectories, (3.13) is evaluated for pursuer and evader. Then, the resultant integral is found as

$$U(\mathbf{v}^{i*}) = \lambda^i (\nabla V^*)^T \tanh(\mathbf{v}^{i*}) - 2\lambda^i \underline{\mathbf{R}}^i \log(\cosh(\mathbf{v}^{i*})) \quad \forall i = p, e \quad (3.60)$$

where  $\log(\cdot)$  is the natural logarithm,  $\underline{\mathbf{R}}^i \triangleq \text{diag}(\mathbf{R}^i) = [r_1^i \ r_2^i \ r_3^i]^T$ , and  $\mathbf{v}^{i*}$  stands for the optimal velocity policy given by (3.17) and (3.18)  $\forall i = p, e$ .

When the evader is moving with the sub-optimal velocity, we set  $U(\pi^e(\boldsymbol{\delta}))$  term in (3.11) to zero and thereby we obtain Hamilton Jacobi Bellman (HJB) equation instead of HJI (3.19). Notice that HJB equation in this case stands for the single player game where the pursuer is the only player. Furthermore, the existence of unique Nash equilibrium by Theorem 1 implies that the value functional (3.11) is convex in  $\mathbf{v}^{p*}$  for  $|v_j^e| \leq \lambda^e \ \forall j \in \{1, 2, 3\}$  given in (3.13), and the functional (3.11) is concave in  $\mathbf{v}^{e*}$  for  $|v_j^p| \leq \lambda^p \ \forall j \in \{1, 2, 3\}$ . Then, (3.11) is separable, and solution of

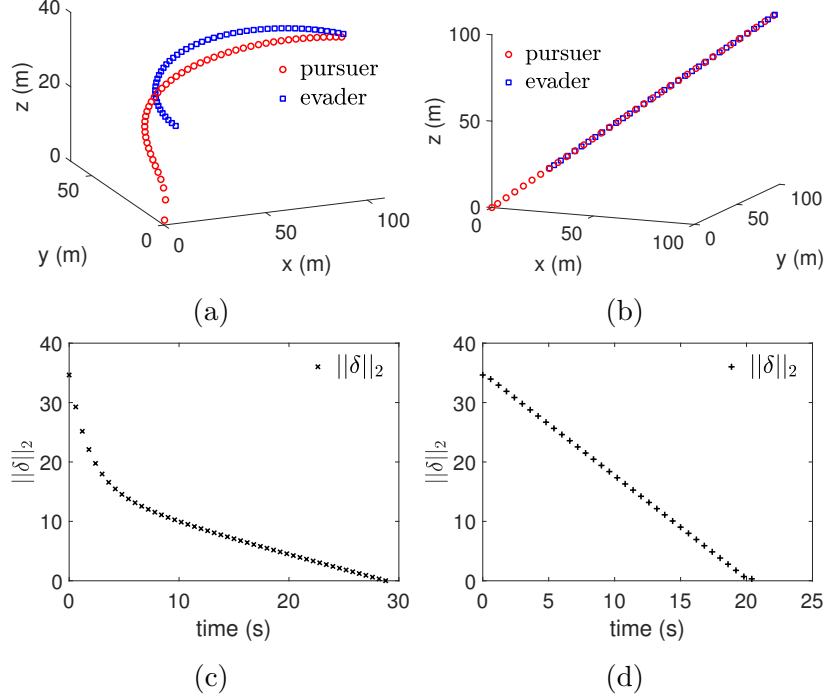


Figure 17: Position of the pursuer and evader: (a)  $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ ,  $\pi^e(\boldsymbol{\delta}) \triangleq \text{suboptimal}$ , (b)  $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ ,  $\pi^e(\boldsymbol{\delta}) \triangleq (3.18)$ .  $L_2$  norm of (3.10): (c)  $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ ,  $\pi^e(\boldsymbol{\delta}) \triangleq \text{suboptimal}$ , (d)  $\pi^p(\boldsymbol{\delta}) \triangleq (3.17)$ ,  $\pi^e(\boldsymbol{\delta}) \triangleq (3.18)$ .

the HJB in terms of the optimal velocity policy for the pursuer remains the same as (3.17).

We conducted two simulation scenarios to validate the proposed methods in this paper. We first consider ZS game with the value functional (3.11), and get the players track desired game optimal velocity trajectories (3.17), (3.18) by selecting ideal forces of players derived in (3.7) and corresponding moments (3.59). Then, we set  $U(\pi^e(\boldsymbol{\delta}))$  term in (3.11) to zero, and by solving the corresponding HJB equation, we played single-player game where the pursuer is the only player. Fig. 27 shows the trajectories followed by the players for each of these scenarios.

In these simulations (Figs. 27 and 28), parameters of the system (3.1), (3.2) are selected as  $m^i = 1kg$ ,  $g = 9.81m/s^2$ ,  $\mathbf{I}_B^i = \mathbf{I}_{3 \times 3}$ , where  $\mathbf{I}_{3 \times 3}$  is a 3x3 identity

matrix. The backstepping gain  $\mathbf{K}^i = 5\mathbf{I}_{3 \times 3}$ . In addition, the bounds (3.13) are  $\lambda^p = 5, \lambda^e = 4$ , and value functional parameters (3.11) are  $\mathbf{Q} = 3\mathbf{I}_{3 \times 3}$ ,  $\mathbf{R}^p = 0.1\mathbf{I}_{3 \times 3}$ ,  $\mathbf{R}^e = 0.125\mathbf{I}_{3 \times 3}$ . The position data of each player is collected through each iteration over the period  $T = 0.01s$ . Lastly  $r^e + r^p$  (3.35) and shown in Fig. 26 is selected as  $0.25m$ .

Notice that Fig. 27 shows the trajectories of the players (Figs. 17a and 17b), and corresponding  $L_2$  norm of the position offset (Figs. 17c and 17d). In addition, regarding the optimal velocity policies for the pursuer and evader, i.e. when  $\pi^p(\boldsymbol{\delta}) \triangleq$  (3.17),  $\pi^e(\boldsymbol{\delta}) \triangleq$  (3.18), Fig. 28 illustrate optimal velocities (3.17),(3.18), control forces (3.47),  $L_2$  norm of velocity error (3.5), and Euler angles (3.48),(3.49).

Fig. 30 shows the simulation results of the PE game when the velocity bounds are  $\lambda^p = 10, \lambda^e = 9$ , and other simulation parameters remain the same as in the PE game illustrated in Figs. 27 and 28.

New Solution for H-infinity Static Output-Feedback Control Using Integral  
Reinforcement Learning.

Yusuf Kartal, Wenqian Xue, Atilla Dogan, Frank Lewis, 2021. *Under Review in  
IEEE Transactions on Cybernetics.*

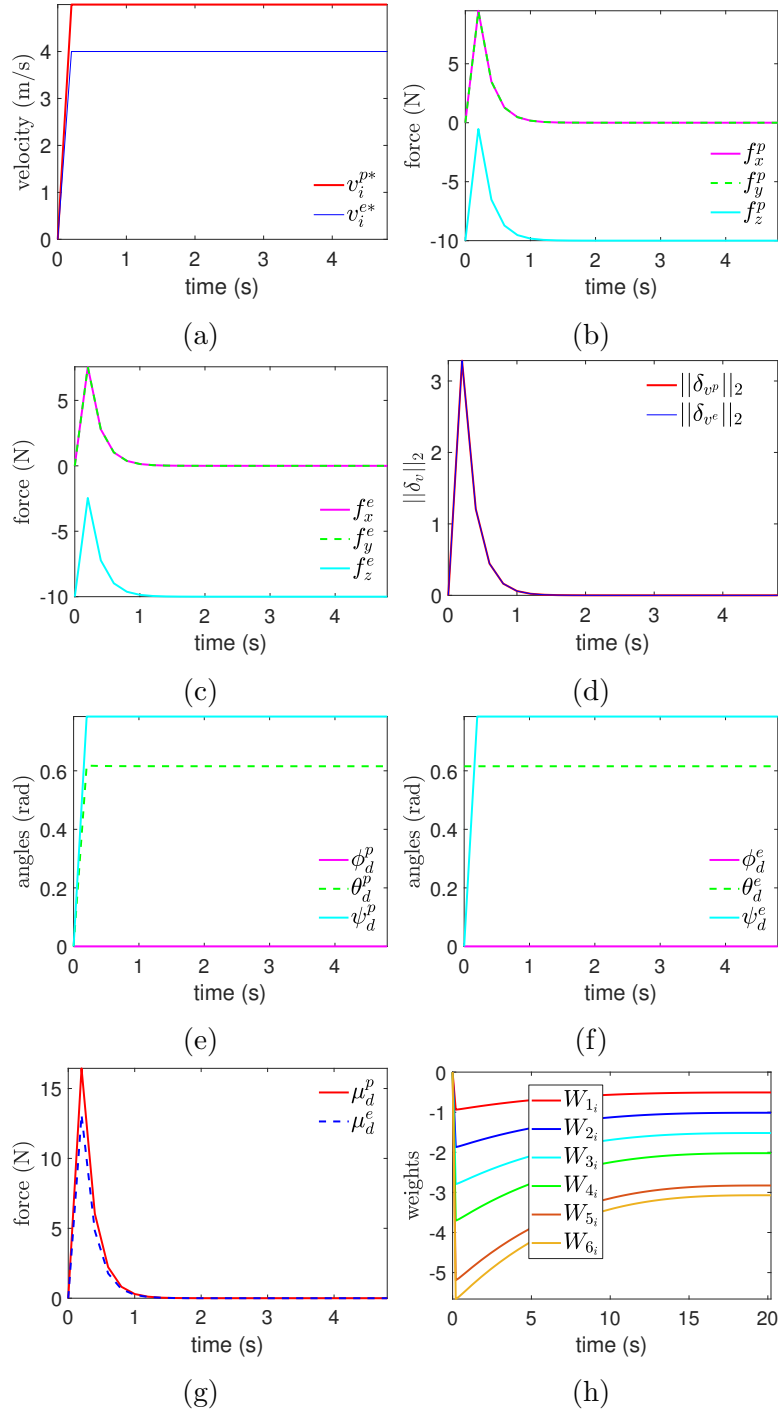


Figure 18: PE game when  $\lambda^p = 5$ ,  $\lambda^e = 4$ : (a) game-optimal velocity of the pursuer and evader  $\forall i \in \{x, y, z\}$ , (b) control force of the pursuer (3.47), (c) control force of the evader (3.47), (d)  $L_2$  norm of (3.5) for each player, (e) Euler angles of the pursuer by (3.47), (f) Euler angles of the evader by (3.47), (g) body evaluated control force of the players (3.50), (h) weights  $\forall i \in \{x, y, z\}$  convergence for the critic NN (3.45).

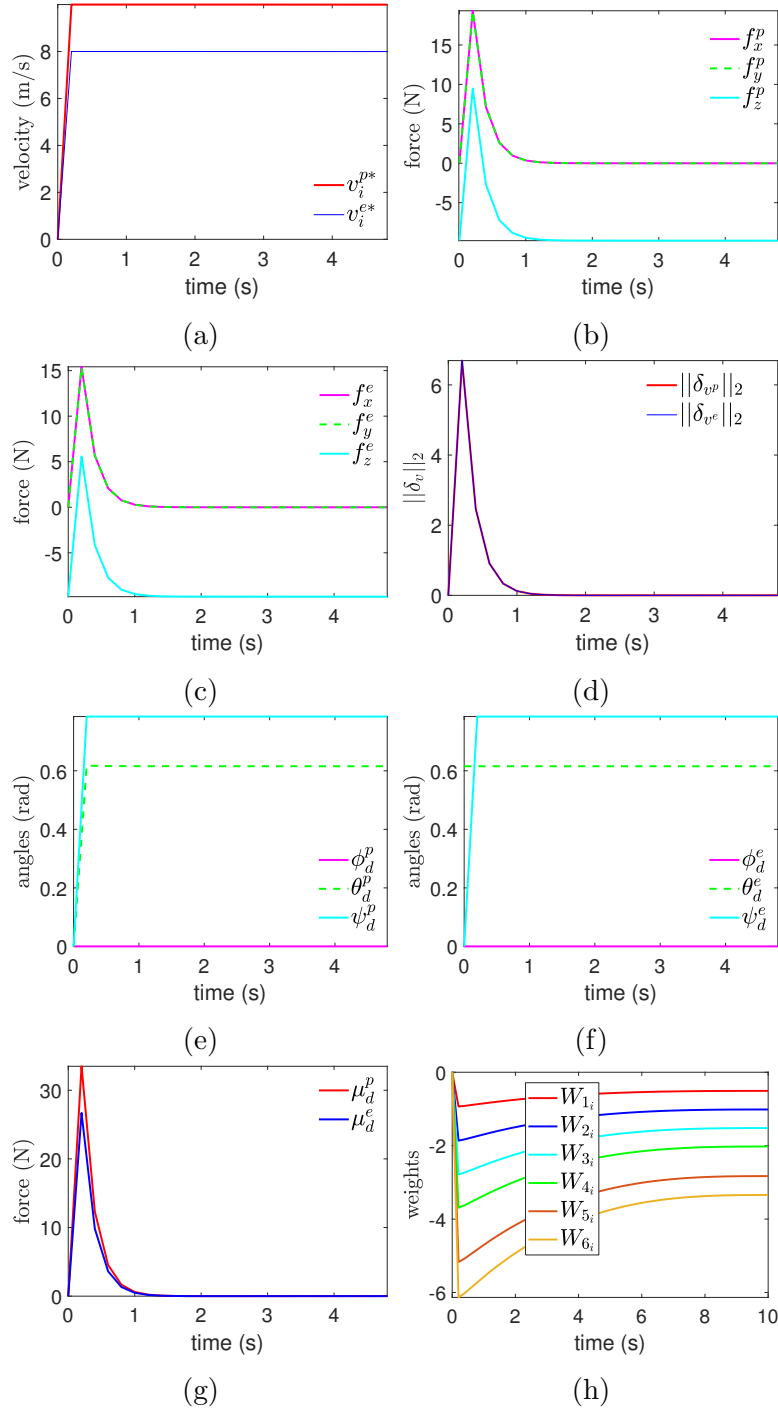


Figure 19: PE game when  $\lambda^p = 10$ ,  $\lambda^e = 8$ : (a) game-optimal velocity of the pursuer and evader  $\forall i \in \{x, y, z\}$ , (b) control force of the pursuer (3.47), (c) control force of the evader (3.47), (d)  $L_2$  norm of (3.5) for each player, (e) Euler angles of the pursuer by (3.47), (f) Euler angles of the evader by (3.47), (g) body evaluated control force of the players (3.50), (h) weights  $\forall i \in \{x, y, z\}$  convergence for the critic NN (3.45).



## CHAPTER 4

### New Solution for H-infinity Static Output-Feedback Control Using Integral Reinforcement Learning

One of the primary objectives in control system design is often to seek a stabilizing controller to regulate the output of a system that experiences disturbances. However, stability is not the only requirement in control system design. An  $L_2$  gain bound of the system, optimality of the control method and detectability of unknown system parameters are other common design specifications. Existing solutions of  $H_\infty$  static output-feedback (OPFB) yield stability and bounded  $L_2$  gain, but employ non-Nash equilibrium solutions. The new formulation of  $H_\infty$  static OPFB control method developed in this paper, is a key to meet these requirements since it guarantees  $L_2$  gain bound of the system by a prescribed attenuation level, asymptotic stability of equilibrium point, and also Nash equilibrium solutions. Then, the integral reinforcement learning technique based on  $H_\infty$  static OPFB Lyapunov iterations is presented to deal with unknown system parameters.

$H_\infty$  control methods has been widely studied in the literature, [10], [11], [12], [13], [14], [15] due to their applicability in variety of engineering areas. Some of these methods guarantee stability and bounded  $L_2$  gain, but an extra condition is required to yield Nash equilibrium. [10] uses this method to design a gain-scheduled normal acceleration control loop for an air-launched unmanned aerial vehicle. Authors of [11] apply this control method on an industrial-type mass spring damper system. The efficacy of control law and the disturbance accommodation properties are shown on a rotor-craft design example in [16]. Moreover, [80] develop an autopilot controller

for an F-16 aircraft by using the  $H_\infty$  static OPFB control method on a linear discrete-time system.

During the last few years, reinforcement learning (RL) algorithms [81], [82] has been used extensively to replace model parameters with a collected system's data [83], [84], [18], [85], [86], [87]. Particularly, offline iterative RL algorithms were studied in [88], [89], and [90]. The work by [88] consider the two-player policy iterations to solve for the feedback strategies of a continuous-time zero-sum game [91] in a sub-optimal manner that requires complete knowledge of system parameters. [84] presented an online RL algorithm to solve the linear quadratic tracking (LQT) problem for partially-unknown continuous-time systems. In [92], authors prove the convergence of IRL algorithm to a sub-optimal OPFB solution without considering the disturbance term when the drift dynamics are unknown. The optimal average cost learning framework is introduced to solve output regulation problem for linear systems with unknown dynamics is studied in [93].

To design an efficient RL algorithm and achieve data-driven optimal control, the RL-based controller designs with neural networks (NNs) in an actor-critic structure [94], critic-only form [79] are proposed. In the off-policy RL algorithm, the system data, which are used to learn the solution of corresponding Hamiltonian, can be generated with arbitrary policies rather than the evaluating policy. This approach is suitably implemented using NNs in [17] and [18].

Standard existing solutions to the static OPFB regulator problem [10], [92], [95], [96], [97] require some additive gain matrix to prove the stability of equilibrium. Unfortunately, this results in the non-Nash equilibrium solutions. In this paper, we propose a novel augmented Hamiltonian, and develop a new iterative algorithms based on stationarity conditions of the augmented Hamiltonian to obtain Nash equilibrium solutions. The salient contributions of this paper summed up into four categories as:

- This paper presents a new solution of the  $H_\infty$  static OPFB control problem. This solution guarantees not only stability and bounded  $L_2$  gain but also Nash equilibrium solutions considering corresponding min-max game.
- Our new solution for static OPFB may have a solution when there is no static OPFB solution using existing techniques.
- Two off-line iterative solution algorithms are given. The first algorithm is based on the Kleinman's algorithm and updates the disturbance gain term. A second algorithm is developed to get the Kleinman's algorithm in the IRL applicable form that only updates the control gain.
- Two off-policy Integral Reinforcement Learning (IRL) algorithms are developed based on the stationarity conditions of an augmented Hamiltonian. This enables the designer to learn the Nash equilibrium solution online without requiring any knowledge of system dynamics' state, control, and disturbance matrices.

The rest of paper is organized as follows. In Section 4.1, preliminaries on control design requirements are introduced, and the formulation of  $H_\infty$  static OPFB control presented. A new solution of optimal  $H_\infty$  control problem and corresponding offline iterative solution algorithms are given in Section 4.2. In Section 4.3, an online off-policy IRL algorithm is developed based on stationarity conditions obtained in Section 4.2. Finally, we have shown the effectiveness of proposed algorithms by applying them to the linearized lateral dynamics of the F-16 aircraft at a particular flight condition in Section 5.

**Notations.** We use the following notations throughout this paper  $\mathbf{I}_n \in \mathbb{R}^{n \times n}$  is the identity matrix. The condition  $A > 0$  ( $\geq 0$ ) denotes the positive (semi) definiteness of a matrix. The operator  $tr()$  denotes trace of a matrix.  $\mathbf{C}^+ = \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}$  is the right-inverse of the full row-rank matrix  $\mathbf{C}$  and the Kronecker product operator is denoted by  $\otimes$ . The determinant of a square matrix is denoted by  $|\cdot|$ .  $vec(A)$  stands

for the  $mn$ -vector formed by stacking the columns of  $A \in \mathbb{R}^{n \times m}$  on top of one another, i.e.,  $\text{vec}(A) = [a_1^T \dots a_m^T]^T$  where  $a_i \in \mathbb{R}^n$  are the columns of  $A$ . Lastly,  $\text{diag}(\zeta_i)$  represents a diagonal matrix with  $\zeta_i \forall i \in 1, \dots, N$  on its diagonal.

#### 4.1 Preliminaries and problem formulation

In this section, preliminaries on Linear Time Invariant (LTI) system, and the corresponding controller design requirements are first introduced. Then, the problem description is presented.

##### 4.1.1 System description and definitions

This section introduces system dynamics and performance specifications that are of interest. Consider the state-space representation of the continuous-time LTI system as

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} + \mathbf{D}\mathbf{d} \\ \mathbf{y} &= \mathbf{C}\mathbf{x}\end{aligned}\tag{4.1}$$

where  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{D} \in \mathbb{R}^{n \times p}$  are system-state, input, disturbance matrices, and  $\mathbf{C} \in \mathbb{R}^{q \times n}$  is assumed to be a full row-rank output matrix to avoid redundant measurements. The corresponding vectors  $\mathbf{x}(t)$ ,  $\mathbf{u}(t)$ ,  $\mathbf{d}(t)$ , and  $\mathbf{y}(t)$  stand for the state, input, disturbance and output respectively.

**Assumption 3.** *The pair  $(\mathbf{A}, \mathbf{B})$  is stabilizable and the pair  $(\mathbf{A}, \mathbf{C})$  is detectable.*

**Assumption 4.** *The system (4.1) is OPFB stabilizable in the sense that the row-space of output matrix  $\mathbf{C}$  contains the sub-space that is spanned by the right eigenvectors correspond to the unstable modes of  $A$ .*

**Assumption 5.** *The non-zero columns of the output matrix  $\mathbf{C}$  are linearly independent.*

**Remark 7.** *The Assumption 7 can be interpreted such that all unstable modes are measured by the output matrix  $\mathbf{C}$  that represents the sensors installed in the system (4.1). The Assumption 8 enables us to recover a state element  $x_i$  precisely from the output vector  $\mathbf{y}$  once it is left multiplied with  $\mathbf{C}^+$ .*

Now, define the fictitious performance output  $\mathbf{z}(t)$  that satisfies

$$\|\mathbf{z}\|_2^2 = \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} \quad (4.2)$$

where  $\mathbf{Q} \geq 0$  and  $\mathbf{R} > 0$  are symmetric design matrices with appropriate dimensions. We assume that  $\mathbf{Q}$  is selected such that the pair  $(\mathbf{A}, \sqrt{\mathbf{Q}})$  is *observable*, which is a standard assumption [60]. Using the property  $\|\mathbf{d}\|_2^2 = \mathbf{d}^T \mathbf{d}$ , a realization of the following inequality  $\forall \mathbf{d} \in [0, \infty)$  implies that the system  $L_2$ -gain is bounded by a prescribed disturbance attenuation level denoted by  $\gamma$

$$\int_0^\infty \|\mathbf{z}\|_2^2 dt \leq \gamma^2 \int_0^\infty \|\mathbf{d}\|_2^2 dt + \beta \quad (4.3)$$

for any non-zero energy-bounded disturbance input  $\mathbf{d}$  [98] where  $\beta$  is a non-negative constant. The condition (5.12) is also called as non-nexpansivity constraint in [74]. Call  $\gamma^*$  the minimum gain for which this occurs. In [99] and [100], an algorithm to find  $\gamma^*$  is given, and a formulation for explicit  $\gamma^*$  that depends on Riccati equation solution is derived for LTI systems under some assumptions. In this paper, we assume that the attenuation level is prescribed and satisfies  $\gamma > \gamma^*$ .

A static OPFB control to regulate the system (4.1) is

$$\mathbf{u} = -\mathbf{K} \mathbf{y} = -\mathbf{K} \mathbf{C} \mathbf{x} \quad (4.4)$$

where  $\mathbf{K} \in \mathbb{R}^{m \times q}$  is the gain matrix. Note that main objective of  $H_\infty$  control using OPFB is to find the stabilizing  $\mathbf{K}$  in an optimal manner while satisfying the condition (5.12), which can be achieved by solving corresponding Hamilton-Jacobi-Isaacs (HJI) equation.

#### 4.1.2 Problem formulation and existing solution of static OPFB Problem

In this section, we relate zero-sum differential game theory to the static OPFB regulation problem in a global optimal manner by revealing various definitions. To satisfy  $L_2$ -gain bound (5.12) with the stabilizing gain in (5.13), an objective functional defined as

$$J(\mathbf{u}, \mathbf{d}) = \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} - \gamma^2 \mathbf{d}^T \mathbf{d}) d\tau. \quad (4.5)$$

Now  $H_\infty$  control problem can be represented as a two-player zero-sum differential game by treating  $\mathbf{u}(t)$  as a minimizing player, whereas  $\mathbf{d}(t)$  maximizing player of (4.5). Then, the game can be formulated as

$$V(\mathbf{x}(0)) = J(\mathbf{u}^*, \mathbf{d}^*) = \min_{\mathbf{u}} \max_{\mathbf{d}} J(\mathbf{u}, \mathbf{d}) \quad (4.6)$$

where  $V(\mathbf{x})$  denotes the value functional corresponding to (4.5) such that

$$V(\mathbf{x}) = \int_t^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} - \gamma^2 \mathbf{d}^T \mathbf{d}) d\tau. \quad (4.7)$$

and the pair  $(\mathbf{u}^*, \mathbf{d}^*)$  denotes the game theoretic saddle point. The game of this kind admits a unique solution pair  $(\mathbf{u}^*, \mathbf{d}^*)$ , if the following Nash condition holds

$$\min_{\mathbf{u}} \max_{\mathbf{d}} J(\mathbf{u}, \mathbf{d}) = \max_{\mathbf{d}} \min_{\mathbf{u}} J(\mathbf{u}, \mathbf{d}). \quad (4.8)$$

The next theorem recalls the necessary and sufficient conditions for the sub-optimal  $H_\infty$  OPFB control method [10].

**Theorem 5.** *The system (4.1) is OPFB stable using the control  $\mathbf{u}_o^e = -\mathbf{K}_o^e \mathbf{y}$  with  $L_2$  gain bounded by  $\gamma > \gamma^*$  if and only if*

1.  $(\mathbf{A}, \mathbf{B})$  is stabilizable and  $(\mathbf{A}, \mathbf{C})$  is detectable.
2. There exists  $\mathbf{K}_o^e$  and  $\mathbf{L}$  such that

$$\mathbf{K}_o^e = \mathbf{R}^{-1}(\mathbf{B}^T \mathbf{P} + \mathbf{L})\mathbf{C}^+ \quad (4.9)$$

where  $\mathbf{P} = \mathbf{P}^T > \mathbf{0}$  is the solution of Riccati equation

$$\mathbf{P}\mathbf{A} + \mathbf{A}^T \mathbf{P} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \mathbf{P} + \mathbf{Q} + \gamma^{-2}\mathbf{P}\mathbf{D}\mathbf{D}^T \mathbf{P} + \mathbf{L}^T \mathbf{R}^{-1} \mathbf{L} = \mathbf{0}. \quad (4.10)$$

**Proof.** See [10] and [98] for the same proof. □

**Remark 8.** *On the one hand, introducing the additive gain matrix  $\mathbf{L}$  provides the extra design freedom. Note that if  $\mathbf{L} = \mathbf{0}$  there may not be a stabilizing solution to (4.9),(5.22) (See Theorem 1 in [10]). On the other hand, this results in a sub-optimal solution for the gain  $\mathbf{K}^s$  (4.9). The Nash equilibrium gain occurs if the Theorem 5 holds with  $\mathbf{L} = \mathbf{0}$ .*

The following lemma recalls the Nash equilibrium solution for the standard  $H_\infty$  control problem.

**Lemma 1.** *The pair  $(u^*, d^*)$  constitutes the Nash equilibrium of the game (4.8) such that*

$$\mathbf{u}^* = -\mathbf{K}^* \mathbf{x} = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} \mathbf{x}, \quad (4.11)$$

$$\mathbf{d}^* = \gamma^{-2} \mathbf{D}^T \mathbf{P} \mathbf{x}. \quad (4.12)$$

**Proof.** Begin with deriving the Hamiltonian to solve the game theoretic saddle point or Nash equilibrium strategy of the game (4.8) as

$$H(\mathbf{V}_x, \mathbf{u}, \mathbf{d}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} - \gamma^2 \mathbf{d}^T \mathbf{d} + \mathbf{V}_x (\mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u} + \mathbf{D} \mathbf{d}) \quad (4.13)$$

where  $\mathbf{V}_x = \partial V / \partial \mathbf{x}$  is the co-state vector. Using the quadratic form  $V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x}$ , and applying the stationarity conditions  $\partial H(\mathbf{V}_x, \mathbf{u}, \mathbf{d}) / \partial \mathbf{u} = \mathbf{0}$  and  $\partial H(\mathbf{V}_x, \mathbf{u}, \mathbf{d}) / \partial \mathbf{d} = \mathbf{0}$  yields the optimal control and disturbance respectively as (5.20) and (5.21).

Notice that the sign of Hessians,  $\partial^2 H(\mathbf{V}_x, \mathbf{u}, \mathbf{d}) / \partial \mathbf{u}^2 > 0$  and  $\partial^2 H(\mathbf{V}_x, \mathbf{u}, \mathbf{d}) / \partial \mathbf{d}^2 < 0$ , along with unboundedness of the limits  $\lim_{d \rightarrow \infty} J(\mathbf{u}^*, \mathbf{d})$ ,  $\lim_{u \rightarrow \infty} J(\mathbf{u}, \mathbf{d}^*)$  indeed show that (5.20) and (5.21) are the global optimal minimizing and maximizing extrema respectively [61]. This indeed verifies that the pair  $(\mathbf{u}^*, \mathbf{d}^*)$  denotes the Nash equilibrium point, which completes the proof.  $\square$

**Remark 9.** *The HJI equation can be obtained as  $H(\mathbf{V}_x, \mathbf{u}^*, \mathbf{d}^*) = 0$  with the boundary condition  $V(0) = 0$ , which also verifies the sub-optimality of gain expression (4.9). Additionally, considering Nash equilibrium policies (5.20) and (5.21), the Game Algebraic Riccati Equation (GARE) is obtained as*

$$\mathbf{P} \mathbf{A} + \mathbf{A}^T \mathbf{P} - \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} + \mathbf{Q} + \gamma^{-2} \mathbf{P} \mathbf{D} \mathbf{D}^T \mathbf{P} = \mathbf{0}. \quad (4.14)$$

If there is no  $\mathbf{K}^s$  to satisfy (4.9) and (5.22), the OPFB  $H_\infty$  control problem may not have even a sub-optimal solution. In the next section, we rigorously analyze this and reveal some novel results.



## 4.2 New solution of OPFB $H_\infty$ game

Notice that Theorem 5 provides necessary and sufficient conditions for static OPFB in a sub-optimal manner. Thence, this does not yield a Nash equilibrium solution unless  $L = 0$ . Moreover, there may not even exist a static OPFB solution to the equations in Theorem 5. In this main section, two methods are proposed to solve  $H_\infty$  OPFB problem. The first method parameterizes the state feedback gains by using the Nash strategies, and applies them to the OPFB design. The second method derives the new optimal  $H_\infty$  OPFB regulator formulation by introducing an augmented Hamiltonian. This method is introduced by [60], but is highly overlooked in the literature. However, it appears to be instrumental in  $H_\infty$  regulator design.

### 4.2.1 Necessary and Sufficient Conditions for the Stabilizing Nash Gain

Herein, we first parameterize static state feedback gains, and then explain how to apply them to the  $H_\infty$  OPFB design.

Note that the Assumptions 7 and 8 are missing in the papers [10], [11], [16], and [98]. However, they are indeed required when applying state feedback gains to the OPFB design.

**Theorem 6.** *Given the necessary conditions in the Assumption 3 and the sufficient condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq \gamma^{-2}\mathbf{D}\mathbf{D}^T$ . The system (4.1) is asymptotically stable using the control  $\mathbf{u}^* = -\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{x}$  (5.20) with  $d = 0$ , and  $L_2$  gain bounded by  $\gamma \forall \|\mathbf{d}\|_2 \in (0, \infty)$ .*

**Proof.** To prove  $L_2$  gain bound condition (5.12), first re-write the Hamiltonian (4.13) by completing the squares as

$$\begin{aligned} H(\mathbf{V}_x, \mathbf{u}, \mathbf{d}) = & H(\mathbf{V}_x, \mathbf{u}^*, \mathbf{d}^*) + (u - u^*)^T R(u - u^*) \\ & - \gamma^{-2}(d - d^*)^T(d - d^*). \end{aligned} \quad (4.15)$$

Then, the objective functional can be re-expressed as

$$\begin{aligned}
J(\mathbf{u}, \mathbf{d}) &= \int_0^\infty \left( H(\mathbf{V}_x, \mathbf{u}^*, \mathbf{d}^*) + (u - u^*)^T R(u - u^*) \right. \\
&\quad \left. - \gamma^{-2}(d - d^*)^T (d - d^*) \right) dt + V(x(0)).
\end{aligned} \tag{4.16}$$

Realize that the non-expansivity constraint (5.12) implies that  $J(u, d) \leq \beta$ . Select  $\beta = V(x(0))$ ,  $u = u^*$ , and note that HJI equation  $H(\mathbf{V}_x, \mathbf{u}^*, \mathbf{d}^*) = 0$  holds with the boundary condition  $V(0)=0$ . Then, (5.52) reduces to

$$\begin{aligned}
J(\mathbf{u}^*, \mathbf{d}) &= - \int_0^\infty \gamma^{-2}(d - d^*)^T (d - d^*) dt + V(x(0)), \\
&\leq V(x(0)), \quad \forall \|\mathbf{d}\|_2 \in (0, \infty).
\end{aligned} \tag{4.17}$$

This proves that the  $L_2$  gain bound condition (5.12) holds with  $\beta = V(x(0))$ ,  $\mathbf{u} = \mathbf{u}^*$ . Additionally, the value of game (4.8) with  $u = u^*$  and  $d = d^*$  is  $V(x(0))$  by (5.52).

To verify asymptotic stability, consider Lyapunov function  $V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x}$  where  $\mathbf{P}$  is solution of the GARE (4.14). Note that for  $\mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \geq \gamma^{-2} \mathbf{D} \mathbf{D}^T$ , the GARE (4.14) has a unique stabilizing solution  $\mathbf{P} = \mathbf{P}^T$ , which is indeed positive definite. This is illustrated in the following realization

$$\begin{aligned}
\dot{V} &= \dot{\mathbf{x}}^T \mathbf{P} \mathbf{x} + \mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} \\
&= \mathbf{x}^T \left( -\mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} - \mathbf{Q} + \gamma^{-2} \mathbf{P} \mathbf{D} \mathbf{D}^T \mathbf{P} \right) \mathbf{x} \\
&\leq -\mathbf{x}^T \mathbf{Q} \mathbf{x} \quad \Leftarrow \quad \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \geq \gamma^{-2} \mathbf{D} \mathbf{D}^T.
\end{aligned} \tag{4.18}$$

Then, the observability of  $(\mathbf{A}, \mathbf{Q})$  verifies that the undisturbed system (4.1), i.e.,  $\mathbf{d} = \mathbf{0}$ , is asymptotically stable by LaSalle's invariance principle [74]. This completes the proof.  $\square$

**Corrolary 1.** *To verify asymptotic stability of the disturbed system, benefit from gain margin  $[c_{\text{lower}}, \infty)$  with  $c_{\text{lower}} < \frac{1}{2}$  property of the  $\mathcal{H}_\infty$  control [74]. Note that the  $\mathcal{H}_\infty$  has gain margin less than  $\frac{1}{2}$  by Chapter 10 in [74] but the lower bound  $c_{\text{lower}}$  is not precisely defined. Then, if the sufficient condition in Theorem 6 is strengthened as  $BR^{-1}B^T \geq 2\gamma^{-2}DD^T$ , the disturbed system (4.1) with  $d = d^*$ , becomes asymptotically stable. Thence, the closed-loop matrix  $\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} + \gamma^{-2}\mathbf{D}\mathbf{D}^T\mathbf{P}$  becomes Hurwitz.*

Till now, we actually parameterize all the stabilizing static state-feedback gains since we used the Nash strategies (5.20) and (5.21) in the proof of Theorem 6. The following Remark explains how to apply them to the  $H_\infty$  OPFB design.

**Remark 10.** *Instead of Nash strategies (5.20) and (5.21), assume that the control and disturbance are selected as*

$$\mathbf{u}_o^* = -(R^{-1}B^T PC^+)C(C^+y) = -\underbrace{(R^{-1}B^T PC^+)}_{K_o^*}y \quad (4.19)$$

$$\mathbf{d}_o^* = (\gamma^{-2}\mathbf{D}^T\mathbf{P}\mathbf{C}^+)C(C^+y). \quad (4.20)$$

*Note that if  $\mathbf{C}$  is an invertible matrix, then all states would be regulated optimally by Theorem 6. On the other hand, if it is not square but full row-rank, then only the states spanned by row space of the output matrix  $\mathbf{C}$  would be regulated optimally given the Assumption 8, and the fact that  $\mathbf{C}^+\mathbf{C}$  projects  $\mathbb{R}^n$  onto the row space of  $\mathbf{C}$ . Additionally, the other states would converge to the origin given the Assumption 7. Realize that the Assumption 3 implies that there could be an unstable mode that is observable but does not belong to the row space of  $\mathbf{C}$ . Therefore, the Assumptions 7 and 8 are indeed required to apply static state feedback gains to the OPFB design. Lastly, the system (4.1) is stable against the worst-case disturbance (4.20), which*

affects only the states that belong to the row space of  $C$  if  $BR^{-1}B^T \geq 2\gamma^{-2}DD^T$  by Corollary 1.

#### 4.2.2 A Direct Method to Obtain OPFB Optimal Gain Solutions

In this section, we propose a new methodology to obtain stabilizing Nash solutions for the  $H_\infty$  OPFB control. This method is direct in the sense that it reaches the same gain solutions as the Section 4.2.1 but do not require two step

Consider the optimal value obtained in the Theorem 6 that corresponds to the zero-effort for each player such that

$$J_0 = \mathbf{x}^T(0)\mathbf{P}\mathbf{x}(0) = tr(\mathbf{P}\mathbf{X}_0) \quad (4.21)$$

where  $tr()$  stands for the trace of a matrix, and  $\mathbf{X}_0 = \mathbf{x}(0)\mathbf{x}^T(0)$ .

The following Lemma is an essential step before introducing the augmented Hamiltonian.

**Lemma 2.** *Given the Assumption 3, let  $\mathbf{K}^a$  be a gain that stabilizes the system (4.1), and the corresponding OPFB control is  $\mathbf{u}_o^a = -\mathbf{K}_o^a\mathbf{y}$ . Additionally, let the disturbance takes the form  $\mathbf{d}_o^a = \mathbf{N}_o\mathbf{y}$  to guarantee it does not affect unobservable modes of the system (4.1). Then the corresponding Lyapunov equation can be derived as*

$$\mathbf{P}\mathbf{A}_c + \mathbf{A}_c^T\mathbf{P} + \mathbf{C}^T\mathbf{K}_o^{aT}\mathbf{R}\mathbf{K}_o^a\mathbf{C} - \gamma^2\mathbf{C}^T\mathbf{N}_o^T\mathbf{N}_o\mathbf{C} + \mathbf{Q} \equiv \mathbf{T} \quad (4.22)$$

where  $\mathbf{T} = \mathbf{T}^T = \mathbf{0}$  and  $\mathbf{A}_c = \mathbf{A} - \mathbf{B}\mathbf{K}_o^a\mathbf{C} + \mathbf{D}\mathbf{N}_o\mathbf{C}$ .

**Proof.** Consider the quadratic form of the value functional (4.5) as  $V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x}$ , and then substitute expressions  $\mathbf{u} = -\mathbf{K}^a \mathbf{y}$  and  $\mathbf{d}_o^a = \mathbf{N}_o \mathbf{y}$  into (4.5) to obtain

$$\mathbf{x}^T \mathbf{P} \mathbf{x} = \int_t^\infty \mathbf{x}^T (\mathbf{C}^T \mathbf{K}_o^{aT} \mathbf{R} \mathbf{K}_o^a \mathbf{C} + \mathbf{Q} - \gamma^2 \mathbf{C}^T \mathbf{N}_o^T \mathbf{N}_o \mathbf{C}) \mathbf{x} d\tau. \quad (4.23)$$

Now, take the derivative of left-side (4.23) and substitute (5.21),  $\mathbf{u} = -\mathbf{K}^a \mathbf{C} \mathbf{x}$  expressions. Lastly, take the derivative of integral in right-side (4.23) using Leibniz's rule, which yields

$$\begin{aligned} \mathbf{x} (\mathbf{A}_c^T \mathbf{P} + \mathbf{P} \mathbf{A}_c) \mathbf{x} = \\ \mathbf{x}^T (-\mathbf{C}^T \mathbf{K}_o^{aT} \mathbf{R} \mathbf{K}_o^a \mathbf{C} - \mathbf{Q} + \gamma^2 \mathbf{C}^T \mathbf{N}_o^T \mathbf{N}_o \mathbf{C}) \mathbf{x}. \end{aligned} \quad (4.24)$$

Realize that the zero equivalent is nothing but the Lyapunov equation given in (4.22). This completes the proof.  $\square$

It is now clear that performing a min-max operation on (4.5), subject to dynamical constraint (4.1) is equivalent to the algebraic problem of finding the pair  $(\mathbf{K}_o^a, \mathbf{N}_o)$  that performs min-max of (4.21) subject to the constraint (4.22). Define following augmented Hamiltonian to solve this modified problem as

$$H^a(\mathbf{K}_o^a, \mathbf{N}_o, \mathbf{S}) = tr(\mathbf{P} \mathbf{X}_0) + tr(\mathbf{T} \mathbf{S}) \quad (4.25)$$

where  $\mathbf{S} \in \mathbb{R}^{n \times n}$  is a symmetric matrix of Lagrange multipliers [60] that needs to be determined, and  $\mathbf{T}$  is given in (4.22).

The next main theorem is a key to find  $\mathbf{S}$  along with a Nash equilibrium control matrix  $\mathbf{K}_o^a$  with respect to (4.25).

**Theorem 7.** *Given the Assumptions 3, 7, and 8, the system (4.1) is asymptotically stable using the control  $\mathbf{u}_o^a = -\mathbf{K}_o^a \mathbf{y}$  with  $d_o^a = 0$  and  $L_2$  gain bounded by  $\gamma$  if*

$$\frac{\partial H^a}{\partial \mathbf{S}} \equiv \mathbf{P} \mathbf{A}_c + \mathbf{A}_c^T \mathbf{P} + \mathbf{C}^T \mathbf{K}_o^{aT} \mathbf{R} \mathbf{K}_o^a \mathbf{C} - \gamma^2 \mathbf{C}^T \mathbf{N}_o^T \mathbf{N}_o \mathbf{C} + \mathbf{Q} = \mathbf{0}, \quad (4.26)$$

$$\frac{\partial H^a}{\partial \mathbf{P}} \equiv \mathbf{S} \mathbf{A}_c^T + \mathbf{A}_c \mathbf{S} + \mathbf{X}_0 = \mathbf{0}, \quad (4.27)$$

$$\frac{\partial H^a}{\partial \mathbf{K}_o^a} \equiv 2\mathbf{R} \mathbf{K}_o^a \mathbf{C} \mathbf{S} \mathbf{C}^T - 2\mathbf{B}^T \mathbf{P} \mathbf{S} \mathbf{C}^T = \mathbf{0}, \quad (4.28)$$

$$\frac{\partial H^a}{\partial \mathbf{N}_o} \equiv -2\gamma^2 \mathbf{N}_o^a \mathbf{C} \mathbf{S} \mathbf{C}^T + 2\mathbf{D}^T \mathbf{P} \mathbf{S} \mathbf{C}^T = \mathbf{0}. \quad (4.29)$$

Furthermore, the following gain expressions solves the  $H_\infty$  static OPFB problem in an optimal manner

$$\mathbf{K}_o^a = \mathbf{K}_o^* = \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} \mathbf{C}^+ \quad (4.30)$$

$$\mathbf{N}_o = \gamma^{-2} \mathbf{D}^T \mathbf{P} \mathbf{C}^+ \quad (4.31)$$

**Proof.** Consider (4.25), which is a constant along the system trajectories since the system (4.1) is LTI and  $\mathbf{z}(t)$  (4.2) is not explicit function of time. This implies that we can apply the constraint test and check the stationarity conditions on the augmented Hamiltonian (4.25) that yields the second-order Lyapunov equation (4.26) and standard Lyapunov equation (4.27) respectively.

Define the variable  $X = \mathbf{x} \mathbf{x}^T$  that includes the system state information, and take the derivative using (4.1) to obtain

$$\mathbf{X} \mathbf{A}_c^T + \mathbf{A}_c \mathbf{X} = \dot{\mathbf{X}}. \quad (4.32)$$

Now, assume that  $A_c$  is Hurwitz. Then, taking the integral of both sides (4.32) from 0 to  $\infty$  yields (4.27) where  $S = \int_0^\infty X dt$ , thereby  $\mathbf{K}_o^a$  should not depend on the solution

$S$ . Therefore, the gain solutions (4.30) and (4.31) are immediate. Realize that the stability of an LTI system does not depend on the initial condition, i.e, local stability implies the global stability. Thence, the gain solutions (4.30) and (4.31) should not depend on  $S$  that depends on the initial condition  $X_0$ . This also verifies our reason to select gain solutions in the given forms, which completes the proof.  $\square$

**Remark 11.** *The Theorem 11 gives the necessary conditions, i.e. the Assumption 3 and 7, and sufficient condition  $\mathbf{BR}^{-1}\mathbf{B}^T \geq \gamma^{-2}\mathbf{DD}^T$  to prove  $L_2$  gain boundedness by a prescribed attenuation level  $\gamma$  (5.12) and OPFB stability considering the worst-case disturbance (4.20). Realize that the condition  $\mathbf{BR}^{-1}\mathbf{B}^T \geq \gamma^{-2}\mathbf{DD}^T$  is only a sufficient condition, which implies that there may be an optimal gain solution which stabilizes (4.1) but does not satisfy  $\mathbf{BR}^{-1}\mathbf{B}^T \geq \gamma^{-2}\mathbf{DD}^T$ . However, in that case, one may not achieve a positive definite solution  $\mathbf{P}$  for the Riccati equation (5.22). Additionally, the positive definiteness of the Riccati equation solution plays a key role for the Kleinman's algorithm in Section 4.2.3, and the IRL algorithm in Section 4.3.*

**Remark 12.** *The Theorem 11 proves that the condition  $\mathbf{BR}^{-1}\mathbf{B}^T \geq \gamma^{-2}\mathbf{DD}^T$  is sufficient to obtain the Nash equilibrium solution. The gain  $\mathbf{K}_o^e$  in Theorem 5 is a sub-optimal stabilizing gain solution with respect to the value functional (4.7). However, the gain solution  $\mathbf{K}^a$  in Theorem 11 always gives a stabilizing Nash equilibrium gain solutions with respect to the game (4.8).*

**Remark 13.** *Note that the system (4.1) is stable in the presence of matched disturbances with the gains (4.30) and (4.31). The unmatched disturbances are only  $L_2$  gain bounded by Theorem 6. Thence, the control (4.19) is robustly stabilizing the equilibrium origin even if the unmatched disturbances exist given Assumptions 7-8.*

The next section proposes an offline model-based algorithm to find the optimal gain solution  $\mathbf{K}_o^a$  iteratively, that plays a key role to develop the IRL Algorithm that will be detailed later in Section 4.3. Note that given the necessary and sufficient

conditions in Theorem 11 and Remark 11, one does not need an iterative solution algorithm to find optimal stabilizing gain  $\mathbf{K}_o^a$  (4.30). However, to develop a model-free algorithm, an iterative solution algorithm is required.

#### 4.2.3 Offline Iterative Solution Algorithms for $H_\infty$ OPFB

This section presents two iterative solution algorithms to obtain minimizing gain (4.30) by using the conditions given in (4.26)-(4.28). In the first algorithm, we employ Kleinman's algorithm (4.26) to obtain Nash equilibrium gain solutions (4.30) and (4.31), whereas in the second algorithm, we use a corresponding Lyapunov equation to not deal with the disturbance gain term  $N_o$ .

The next algorithm performs a sequence of four-step iterations based on the Kleinman's Algorithm [101] to find the optimal control gains (4.30) and (4.31).

**Algorithm 1.** (*Offline iterative solution with Lyapunov equations. Kleinman's Algorithm.*)

1. *Initialize: Set  $k = 1$ ,  $P_0 = 0$ ,  $N_0 = 0$  and given the Assumption 3 and the sufficient condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq 2\gamma^{-2}\mathbf{D}\mathbf{D}^T$ , select a gain  $\mathbf{F}_0$  such that  $\mathbf{A} - \mathbf{B}\mathbf{F}_0$  is asymptotically stable.*
2.  *$k^{\text{th}}$  iteration: Solve for  $\mathbf{P}_k$*

$$\mathbf{0} = \mathbf{P}_k \mathbf{A}_k + \mathbf{A}_k^T \mathbf{P}_k + \mathbf{F}_{k-1}^T \mathbf{R} \mathbf{F}_{k-1} + \mathbf{Q} \quad (4.33)$$

where  $\mathbf{A}_k = \mathbf{A} - \mathbf{B}\mathbf{F}_k + \frac{1}{2}\mathbf{D}\mathbf{N}_k$ . Finally, update the gains

$$\mathbf{F}_k = \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}_k, \quad (4.34)$$

$$\mathbf{N}_k = \gamma^{-2} \mathbf{D}^T \mathbf{P}_k. \quad (4.35)$$



Set the cost  $J_k = \text{tr}(\mathbf{P}_k \mathbf{X}_0)$ .

3. Check: If  $\mathbf{F}_{k-1}$  and  $\mathbf{F}_k$  are close enough to each other, go to step (4) else go to step (2).
4. Terminate: Given the Assumptions 7 and 8, set the OPFB gains  $\mathbf{K}_o^a = \mathbf{F}_k \mathbf{C}^+$ ,  $\mathbf{N}_o = \mathbf{N}_k \mathbf{C}^+$  and the cost  $\mathbf{J} = \mathbf{J}_k$ .  $\square$

Note that the closed-loop stability, and  $L_2$  gain boundedness implies that (4.36) has a unique stabilizing optimal solution  $\mathbf{P} > \mathbf{0}$  by Theorem 6. A comprehensive study for the solution of generalized Riccati equations can be found in [102]. The Algorithm 2 is based on the iterative solution algorithm presented in [103] whose convergence is proved by establishing the connection between Newton's method in [104]. Additionally, by considering the condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq \gamma^{-2}\mathbf{D}\mathbf{D}^T$ , the monotonic convergence, i.e,  $P_k < P_{k-1}$ , is straight-forward from the Theorem 13.5.8 in [104]. A related algorithm is also examined in Chapter 8 of the book [60].

Now, to develop an algorithm that accounts only the optimal gain  $\mathbf{K}_o^a$ , we manipulate the steps of the Algorithm 1. The resultant Algorithm 2 finds the the optimal control gain  $\mathbf{K}_o^a$  (4.30) iteratively.

**Algorithm 2.** (*Offline iterative solution with Lyapunov equations. Modified Kleinman's Algorithm.*)

1. Initialize: Set  $k = 1$ ,  $\mathbf{P}_0 = \mathbf{0}$ , and given the Assumption 3 and the sufficient condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq 2\gamma^{-2}\mathbf{D}\mathbf{D}^T$ , select a gain  $\mathbf{F}_0$  such that  $\mathbf{A} - \mathbf{B}\mathbf{F}_0$  is asymptotically stable.
2.  $k^{\text{th}}$  iteration: Solve for  $\mathbf{P}_k$

$$\mathbf{0} = \mathbf{P}_k \mathbf{A}_k + \mathbf{A}_k^T \mathbf{P}_k + \mathbf{F}_{k-1}^T \mathbf{R} \mathbf{F}_{k-1} + \mathbf{Q} + \gamma^{-2} \mathbf{P}_{k-1} \mathbf{D} \mathbf{D}^T \mathbf{P}_{k-1} \quad (4.36)$$

where  $\mathbf{A}_k = \mathbf{A} - \mathbf{B}\mathbf{F}_k$ . Finally, update the gain

$$\mathbf{F}_k = \mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}_k. \quad (4.37)$$

Set the cost  $J_k = \text{tr}(\mathbf{P}_k\mathbf{X}_0)$ .

3. Check: If  $\mathbf{F}_{k-1}$  and  $\mathbf{F}_k$  are close enough to each other, go to step (4) else go to step (2).
4. Terminate: Given the Assumptions 7 and 8, set the OPFB gain  $\mathbf{K}_0^a = \mathbf{F}_k\mathbf{C}^+$  and the cost  $J = J_k$ . □

The next section uses the Algorithm 2 to develop model-free algorithms considering different scenarios for the availability of system data.

### 4.3 Online integral reinforcement learning solution algorithm for $H_\infty$ OPFB

In this section, we first develop an online off-policy integral reinforcement learning (IRL) algorithm [105], which is a model-free version of the Algorithm 2. This algorithm assumes that the system state data is available to deal with unknown  $A$ ,  $B$  and  $D$  matrices. Then, we develop a novel IRL algorithm that solves the optimal  $H_\infty$  regulator problem by learning the Nash equilibrium gain solution (4.30) without requiring knowledge of the system state data.

#### 4.3.1 Online off-policy IRL algorithms

This section introduces an off-policy IRL algorithm to deal with the unknown system matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{D}$ . In this case, both of Algorithm 1 and 2 lose their applicability as they are model-based. However, we still benefit from the convergence

properties of Algorithm 2 while developing off-policy IRL method. To this end, we represent system dynamics (4.1) as

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}_k + \mathbf{D}\mathbf{d} + \mathbf{B}(\mathbf{u} - \mathbf{u}_k) \\ &= \mathbf{A}_k\mathbf{x} + \mathbf{D}\mathbf{d} + \mathbf{B}(\mathbf{u} - \mathbf{u}_k)\end{aligned}\tag{4.38}$$

where  $\mathbf{u}_k = -\mathbf{F}_k\mathbf{x} \in \mathbb{R}^m$  is the control policy to be updated with  $\mathbf{F}_k$  given in Algorithm 2.

Firstly, to obtain  $\mathbf{P}_k$  without information of the system matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{D}$ , take the derivative of value functional  $V(\mathbf{x}(t)) = \mathbf{x}^T(t)\mathbf{P}\mathbf{x}(t)$  by using the new representation of the system dynamics (4.38)

$$\begin{aligned}\dot{V} &= \mathbf{x}^T \mathbf{A}_k^T \mathbf{P}_k \mathbf{x} + \mathbf{x}^T \mathbf{P}_k \mathbf{A}_k \mathbf{x} \\ &\quad + 2\mathbf{d}^T \mathbf{D}^T \mathbf{P}_k \mathbf{x} + 2(\mathbf{u} + \mathbf{F}_k \mathbf{x})^T \mathbf{B}^T \mathbf{P}_k \mathbf{x}.\end{aligned}\tag{4.39}$$

To employ the approach given in Algorithm [101], we define the following two new variables

$$\mathbf{G}_{k+1} = \gamma^{-2} \mathbf{D}^T \mathbf{P}_k, \quad \mathbf{F}_{k+1} = \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}_k.\tag{4.40}$$

Re-write (4.36) in terms of new variable (4.40) to get the Algorithm 2 in the Kleinman's form as

$$\bar{\mathbf{Q}} = \mathbf{P}_k \mathbf{A}_k + \mathbf{A}_k^T \mathbf{P}_k\tag{4.41}$$

where  $\bar{\mathbf{Q}} = -\mathbf{F}_k^T \mathbf{R} \mathbf{F}_k - \mathbf{Q} - \gamma^2 \mathbf{G}_k^T \mathbf{G}_k$ . Additionally, express (4.39) in terms of the new variables introduced in (4.40) and (4.41) as

$$\dot{V} = 2(\gamma^2 \mathbf{d}^T \mathbf{G}_{k+1} + (\mathbf{u} + \mathbf{F}_k \mathbf{x})^T \mathbf{R} \mathbf{F}_{k+1}) \mathbf{x} + \mathbf{x}^T \bar{\mathbf{Q}} \mathbf{x}.\tag{4.42}$$

Then, integrate both sides from  $t$  to  $t + T$  to obtain

$$\begin{aligned} V(t+T) - V(t) &= \int_t^{t+T} 2(\mathbf{u} + \mathbf{F}_k \mathbf{x})^T \mathbf{R} \mathbf{F}_{k+1} \mathbf{x} d\tau \\ &+ \int_t^{t+T} 2\gamma^2 \mathbf{d}^T \mathbf{G}_{k+1} \mathbf{x} d\tau + \int_t^{t+T} \mathbf{x}^T \bar{\mathbf{Q}} \mathbf{x} d\tau. \end{aligned} \quad (4.43)$$

Based on these manipulations, the online off-policy IRL algorithm can be developed.

Note that this new IRL algorithm and the Algorithm 2 are equivalent. However, on the contrary offline Algorithm 2, the IRL Algorithm 3 does not require information of  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{D}$  matrices. It only requires the sufficient amount of data that belongs  $\mathbf{x}(t)$ ,  $\mathbf{u}(t)$ , and  $\mathbf{d}(t)$  vectors, which can be collected online.

**Algorithm 3.** (*Off-policy IRL algorithm assuming  $x$  is given.*)

1. Initialize: Set  $k = 1$ ,  $\mathbf{G}_0 = \mathbf{0}$  and given Assumption 3 determine a stabilizing gain  $\mathbf{F}_0$ .
2.  $k^{\text{th}}$  iteration: Use (4.43) to update  $\mathbf{P}_k$ ,  $\mathbf{G}_{k+1}$  and  $\mathbf{F}_{k+1}$  simultaneously

$$\begin{aligned} \mathbf{x}^T(t+T) \mathbf{P}_k \mathbf{x}(t+T) - \mathbf{x}^T(t) \mathbf{P}_k \mathbf{x}(t) - \int_t^{t+T} 2\gamma^2 \mathbf{d}^T \mathbf{G}_{k+1} \mathbf{x} d\tau \\ - \int_t^{t+T} 2(\mathbf{u} + \mathbf{F}_k \mathbf{x})^T \mathbf{R} \mathbf{F}_{k+1} \mathbf{x} d\tau = \int_t^{t+T} \mathbf{x}^T \bar{\mathbf{Q}} \mathbf{x} d\tau. \end{aligned} \quad (4.44)$$

Set the cost  $J_k = \text{tr}(\mathbf{P}_k \mathbf{X}_0)$ .

3. Check: If  $\mathbf{F}_k$  and  $\mathbf{F}_{k+1}$  are close enough to each other, go to step (4) else go to step (2).
4. Terminate: Given the Assumptions 7 and 8, set  $\mathbf{K}_o^a = \mathbf{F}_{k+1} \mathbf{C}^+$  and  $\mathbf{J} = \mathbf{J}_k$ .

□

**Remark 14.** Note that right-side of the equation (4.44) in the Algorithm 3 consists of known terms, and hence it can be solved for  $\mathbf{P}_k$ ,  $\mathbf{G}_{k+1}$  and  $\mathbf{F}_{k+1}$  matrices using well

established least-squares technique by converting them to the set of linear equations [17]. Since (4.44) does not use any system matrix information except from  $\mathbf{C}$ , the Algorithm 3 is said to be model-free. Realize that the output matrix  $\mathbf{C}$  represents the sensors placed to the system (4.1), which is clearly known.

Realize that the Algorithm 3 achieves the static OPFB gain assuming  $x$  is available. However, we are only given the output data  $\mathbf{y}$ . Therefore, we need to come up with an another algorithm that does not require system state data. One approach is to make use of the observability matrix while developing a model free algorithm [106]. However, this approach indeed requires the pair  $(\mathbf{A}, \mathbf{C})$  to be observable. Thence, from now on we assume that the detectability condition in the Assumption 3 is strengthened as observability condition.

**Theorem 8.** *For any given  $n$ -dimensional observable system, there exists a sufficiently small time interval  $[0 T]$  such that if  $N$  sampling times satisfy  $0 \leq t - i\Delta t \leq T \forall i \in 1, \dots, N$  where  $\Delta t$  is the delayed time and assumed fixed, then the system is  $N$ -sample observable.*

**Proof.:** See [107] for the same proof.

To make use of the observability matrix define

$$\mathbf{x}(t)^T \mathbf{P} \mathbf{x}(t) = \mathbf{Y}^T(t) \tilde{\mathbf{P}} \mathbf{Y}(t) \quad (4.45)$$

where  $\mathbf{Y}(t) = [\mathbf{y}^T(t) \dot{\mathbf{y}}^T(t) \dots \mathbf{y}^{n-1T}(t)]^T = \mathcal{O} \mathbf{x}$ , and hence  $\tilde{\mathbf{P}} = \mathcal{O}_L^T \mathbf{P} \mathcal{O}_L$  with  $\mathcal{O}_L \in \mathbb{R}^{n \times nq}$  is the left-inverse of the observability matrix  $\mathcal{O} \in \mathbb{R}^{nq \times n}$ . Note that each derivative of the output  $\mathbf{y}$  can be obtained by making use of the Taylor Series expansion on the collected data  $\mathbf{y}(t + i\Delta t)$  around  $\mathbf{y}(t)$ . An example is  $\ddot{\mathbf{y}} = \frac{\mathbf{y}(t+\Delta t) + \mathbf{y}(t-\Delta t) - 2\mathbf{y}(t)}{\Delta t^2}$  where  $\mathbf{y}(t + \Delta t) = \mathbf{y}(t) + \Delta t \dot{\mathbf{y}}(t) + 0.5\Delta^2 t \ddot{\mathbf{y}}(t)$  and  $\mathbf{y}(t - \Delta t) = \mathbf{y}(t) - \Delta t \dot{\mathbf{y}}(t) + 0.5\Delta^2 t \ddot{\mathbf{y}}(t)$ .

Now, select  $Q = kC^T C$  with  $k > 0$  is a scalar, and define the following variables

$$\tilde{\mathbf{F}}_k = \mathbf{F}_k \mathcal{O}_L, \quad \tilde{\mathbf{G}}_k = \mathbf{G}_k \mathcal{O}_L, \quad (4.46)$$

and the known term  $\tilde{\mathbf{Q}} = -\tilde{\mathbf{F}}_k^T \mathbf{R} \tilde{\mathbf{F}}_k - \hat{\mathbf{Q}} - \gamma^2 \tilde{\mathbf{G}}_k^T \tilde{\mathbf{G}}_k$  where  $\hat{\mathbf{Q}} = k[\mathbf{I}_q \ 0 \ \dots \ 0]_{qn \times q}^T [\mathbf{I}_q \ 0 \ \dots \ 0]_{q \times qn}$ .

Realize that  $\mathbf{x}^T \mathbf{Q} \mathbf{x} = \mathbf{y}^T \mathbf{y} = \mathbf{Y}^T \hat{\mathbf{Q}} \mathbf{Y}$  holds when  $\mathbf{Q} = k\mathbf{C}^T \mathbf{C}$ , and it enables us to treat  $\tilde{\mathbf{Q}}$  as a known term. Further, since the pair  $(\mathbf{A}, \mathbf{C})$  is observable  $(\mathbf{A}, k\mathbf{C}^T \mathbf{C})$  is also observable. Based on these manipulations, a new Algorithm 4 can be developed.

**Algorithm 4.** (*Off-policy IRL algorithm,  $\mathbf{x}$  is not required but the pair  $(\mathbf{A}, \mathbf{C})$  must be observable.*)

1. *Initialize: Set  $k = 1$ ,  $\tilde{\mathbf{G}}_0 = \mathbf{0}$  and given Assumption 3 determine a stabilizing gain  $\tilde{\mathbf{F}}_0$ .*
2.  *$k^{\text{th}}$  iteration: Update  $\tilde{\mathbf{P}}_k$ ,  $\tilde{\mathbf{G}}_{k+1}$  and  $\tilde{\mathbf{F}}_{k+1}$  simultaneously*

$$\begin{aligned} & \mathbf{Y}^T(t+T) \tilde{\mathbf{P}}_k \mathbf{Y}(t+T) - \mathbf{Y}^T(t) \tilde{\mathbf{P}}_k \mathbf{Y}(t) - \int_t^{t+T} 2\gamma^2 \mathbf{d}^T \tilde{\mathbf{G}}_{k+1} \mathbf{Y} \, d\tau \\ & - \int_t^{t+T} 2(\mathbf{u} + \tilde{\mathbf{F}}_k \mathbf{Y})^T \mathbf{R} \tilde{\mathbf{F}}_{k+1} \mathbf{Y} \, d\tau = \int_t^{t+T} \mathbf{Y}^T \tilde{\mathbf{Q}} \mathbf{Y} \, d\tau. \end{aligned} \quad (4.47)$$

Set the cost  $J_k = \text{tr}(\tilde{\mathbf{P}}_k \mathbf{X}_0)$ .

3. *Check: If  $\tilde{\mathbf{F}}_k$  and  $\tilde{\mathbf{F}}_{k+1}$  are close enough to each other, go to step (4) else go to step (2).*
4. *Terminate: Given the Assumptions 7 and 8, set  $\mathbf{u} = -\tilde{\mathbf{F}}_{k+1} \mathbf{Y}$  and  $\mathbf{J} = \mathbf{J}_k$ .  $\square$*

Realize that the Algorithm 4 converges with the static state feedback expression since  $\mathbf{u} = -\tilde{\mathbf{F}}_{k+1} \mathbf{Y} = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}_k \mathcal{O}_L \mathcal{O} \mathbf{x} = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}_k \mathbf{x}$  (5.20). Thence, it regulates not only the state elements spanned in the row-space of  $\mathbf{C}$  but all states. However, it creates additional complexity as it requires more data to be collected to converge. On the other hand, the Algorithm 4 does not require system state data

$\mathbf{x}$ , and it directly achieves the Nash equilibrium control  $\mathbf{u}$  by making use of the new variables  $\tilde{\mathbf{F}}_{k+1}$  and  $\mathbf{Y}$  instead of calculating the OPFB gain  $\mathbf{K}_k^a$ . Additionally, since the Algorithm 4 is obtained with the change of variables in the Algorithm 3, it shares the similar convergence properties with the Algorithm 3, and hence the Algorithm 2. The next section shows a way of solving coupled linear equations.

#### 4.3.2 Data-based implementation of the IRL Algorithm

This section introduces a least-squares method to solve step (2) in Algorithm 3. Although, value function approximation is a popular tool and can be employed to solve step 2 in Algorithm 3, it requires three Neural Networks (NNs), i.e., the actor NN to approximate the value functional (4.7), the critic NN to update control policy and the disturber NN to update disturbance policy [18]. This causes a complicated NN design procedure. Instead, we use a least-squares method to solve for  $\mathbf{P}_k$  in step (2) of Algorithm 3, however, we first need to find  $\mathbf{P}_k$ .

Now, we use the Kronecker product property  $\mathbf{a}^T \mathbf{W} \mathbf{b} = (\mathbf{b}^T \otimes \mathbf{a}^T) \text{vec}(\mathbf{W})$  to rewrite (4.44) as

$$\begin{aligned} & [\hat{\mathbf{x}}(t+T) - \hat{\mathbf{x}}(t)]^T \hat{\mathbf{P}}_k - 2\gamma^2 \left( \int_t^{t+T} \mathbf{x}^T \otimes \mathbf{d}^T d\tau \right) \text{vec}(\mathbf{G}_{k+1}) \\ & - 2 \left( \int_t^{t+T} \mathbf{x}^T \otimes [(\mathbf{u} + \mathbf{F}_k \mathbf{x})^T \mathbf{R}] d\tau \right) \text{vec}(\mathbf{F}_{k+1}) = \int_t^{t+T} \mathbf{x}^T \bar{\mathbf{Q}} \mathbf{x} d\tau \end{aligned} \quad (4.48)$$

where the vectors  $\hat{\mathbf{x}} \in \mathbb{R}^{\frac{n(n-1)}{2}}$  and  $\hat{\mathbf{P}}_k \in \mathbb{R}^{\frac{n(n-1)}{2}}$  are defined in the following forms

$$\begin{aligned} \hat{\mathbf{x}} &= [x_1^2, 2x_1x_2, \dots, 2x_1x_n, x_2^2, \dots, 2x_{n-1}x_n, x_n^2]^T \\ \hat{\mathbf{P}}_k &= [P_{k(11)}, \dots, P_{k(1n)}, P_{k(22)}, \dots, P_{k(2n)}, \dots, P_{k(nn)}]^T. \end{aligned} \quad (4.49)$$

To solve  $\mathbf{P}_k$ ,  $\mathbf{G}_{k+1}$  and  $\mathbf{F}_{k+1}$  in (4.48), define

$$\mathbf{d}_x = [\hat{\mathbf{x}}(t+T) - \hat{\mathbf{x}}(t), \dots, \quad (4.50)$$

$$\hat{\mathbf{x}}(t+s_1T) - \hat{\mathbf{x}}(t+(s_1-1)T)]^T \in \mathbb{R}^{s_1 \times \frac{n(n-1)}{2}};$$

$$\mathbf{I}_{xd} = \left[ \int_t^{t+T} (\mathbf{x} \otimes \mathbf{d}) d\tau, \dots, \quad (4.51)$$

$$\int_{t+(s_1-1)T}^{t+s_1T} (\mathbf{x} \otimes \mathbf{d}) d\tau \right]^T \in \mathbb{R}^{s_1 \times np};$$

$$\mathbf{I}_{xu} = \left[ \int_t^{t+T} \mathbf{x} \otimes (\mathbf{u} + \mathbf{F}_k \mathbf{x}) d\tau, \dots, \quad (4.52)$$

$$\int_{t+(s_1-1)T}^{t+s_1T} \mathbf{x} \otimes (\mathbf{u} + \mathbf{F}_k \mathbf{x}) d\tau \right]^T \in \mathbb{R}^{s_1 \times nm};$$

$$\Phi = [\mathbf{d}_x, -2\mathbf{I}_{xd}, -2\mathbf{I}_{xu}(\mathbf{I}_n \otimes \mathbf{R})]; \quad (4.53)$$

$$\Psi = -\left[ \int_t^{t+T} \mathbf{x}^T \bar{\mathbf{Q}} \mathbf{x} d\tau, \dots, \int_{t+(s_1-1)T}^{t+s_1T} \mathbf{x}^T \bar{\mathbf{Q}} \mathbf{x} d\tau \right]^T, \quad (4.54)$$

where integer  $s_1 > 0$  is sampling data group number. Then solution can be obtained by

$$\begin{bmatrix} \hat{\mathbf{P}}_k \\ \text{vec}(\mathbf{G}_{k+1}) \\ \text{vec}(\mathbf{F}_{k+1}) \end{bmatrix} = (\Phi^T \Phi)^{-1} \Phi^T \Psi. \quad (4.55)$$

**Remark 15.** To ensure that (4.55) achieves the unique solution, the persistence of excitation condition needs to be satisfied. To this end, probing noise should be injected to control input  $u$  in (37). This is also called as exploration noise that does not affect the convergence [17]. In addition, the data group number  $s_1$  should be no less than  $\frac{n(n+1)}{2} + np + nm$ , which is the number of unknown parameters to be calculated by (4.55).



In the next section, the correct performance of the proposed methods are validated by applying them to the lateral control control of linearized F-16 dynamics.

#### 4.4 Simulation results

In this section, an example is given to verify the correct performance of proposed algorithms that solves optimal  $H_\infty$  static OPFB regulator problem. We used the F-16 linearized lateral dynamics at a particular flight condition given in example 5.3-1 of the book [108]. Parameters of the linearized F-16 lateral dynamics and the corresponding system vectors are

$$\begin{aligned}
 \mathbf{A} &= \begin{bmatrix} -0.32 & 0.064 & 0.036 & -0.992 & 0 & 0.001 \\ 0 & 0 & 1 & 0.004 & 0 & 0 \\ -30.65 & 0 & -3.68 & 0.665 & -0.733 & 0.132 \\ 8.54 & 0 & -0.025 & -0.476 & -0.032 & -0.062 \\ 0 & 0 & 0 & 0 & -20.2 & 0 \\ 0 & 0 & 0 & 0 & 0 & -20.2 \end{bmatrix}; \\
 \mathbf{B} &= \begin{bmatrix} 0 & 0 & 0 & 0 & 20.2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 20.2 \end{bmatrix}^T; \\
 \mathbf{C} &= \begin{bmatrix} 0 & 0 & 0 & 57.2958 & 0 & 0 \\ 0 & 0 & 57.2958 & 0 & 0 & 0 \\ 57.2958 & 0 & 0 & 0 & 0 & 0 \\ 0 & 57.2958 & 0 & 0 & 0 & 0 \end{bmatrix}; \\
 \mathbf{x} &= \begin{bmatrix} \beta & \phi & p & r & \delta_a & \delta_r \end{bmatrix}^T; \quad \mathbf{u} = \begin{bmatrix} u_a & u_r \end{bmatrix}^T; \\
 \mathbf{y} &= \begin{bmatrix} r & p & \beta & \phi \end{bmatrix}^T
 \end{aligned} \tag{4.56}$$

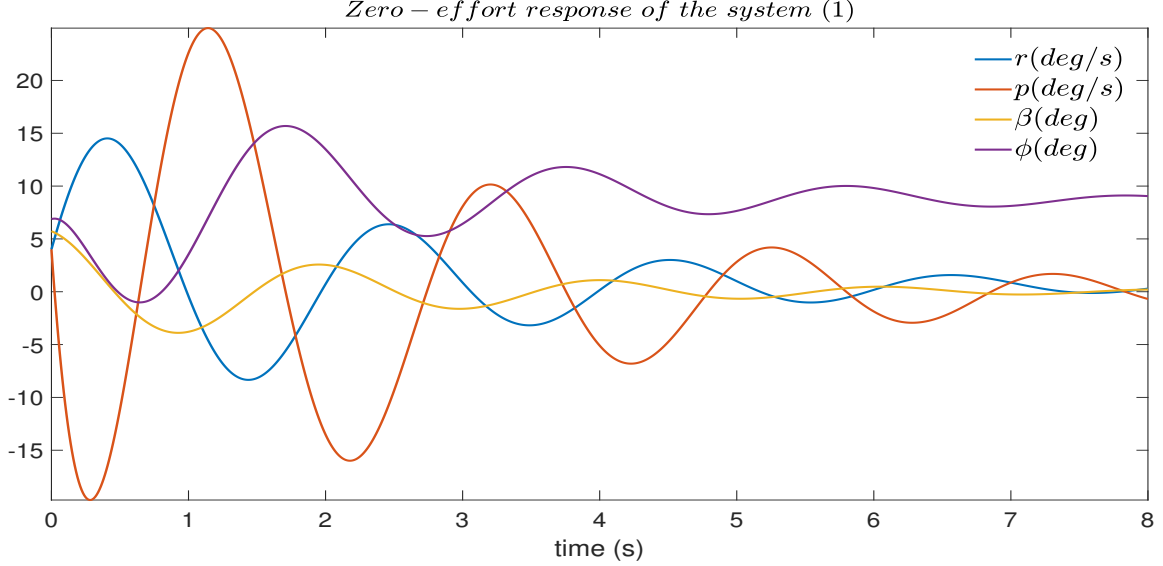


Figure 20: System response when no control is applied.

where  $\beta$  denotes the side-slip,  $\phi$  is the bank angle,  $p$  is the roll rate,  $r$  is the yaw rate,  $\delta_a$  is the aileron actuator angle,  $\delta_r$  is the rudder actuator angle,  $u_a$  is the aileron servo input,  $u_r$  is the rudder servo input. The factor of 57.2958 in the output-matrix  $\mathbf{C}$  converts radians to degrees.

Additionally, the objective functional (4.5) parameters are selected as  $\mathbf{Q} = \text{diag}([50 \ 100 \ 100 \ 50 \ 0 \ 0])$ ,  $\mathbf{R} = \rho \times \text{diag}([0.1 \ 0.7])$  with  $\rho = 0.3$  for computation of the OPFB gain. Also, select the disturbance matrix as

$$\mathbf{D} = \begin{bmatrix} 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{bmatrix}^T, \quad (4.57)$$

and set the attenuation level  $\gamma = 2.5$ . To examine the robustness, assume that the system (4.1) experiences the worst-case disturbance (4.20). Realize that all of the Assumptions 3-8, and the sufficient condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq 2\gamma^{-2}\mathbf{D}\mathbf{D}^T$  are satisfied with the given parameters in (4.56) and (4.57). Now, we first illustrate the

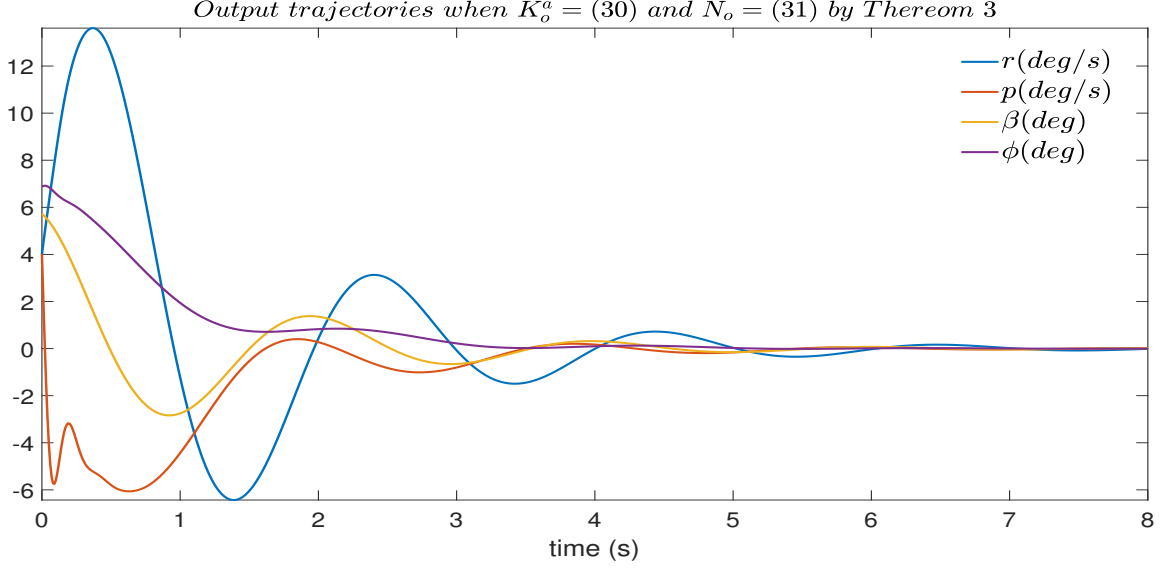


Figure 21: System response when the Nash gains (4.30) and (4.31) are employed.

performance of model-based Nash gain solutions (4.30) and (4.31), and then check whether the the Algorithms 1 and 2 converges the same game solutions as (4.30) and (4.31). After verifying the correctness of them, we illustrate the performance of the Nash solutions obtained in the Algorithms 3 and 4 that are model-free.

The Fig. 26 illustrates the zero control-effort response of the system (4.1). The Fig. 27 illustrates the optimal stabilizing gain (4.30) performance, which is derived in Theorem 11. The Nash OPFB gains found by (4.30) and (4.31) are

$$\mathbf{K}_o^a = \begin{bmatrix} -0.1395 & -0.8714 & 1.4785 & -1.0000 \\ -0.1410 & 0.0273 & -0.1070 & 0.0330 \end{bmatrix}, \quad (4.58)$$

$$\mathbf{N}^o = \begin{bmatrix} -0.0001 & -0.0006 & 0.0011 & -0.0007 \\ -0.0005 & 0.0001 & -0.0004 & 0.0001 \end{bmatrix}. \quad (4.59)$$

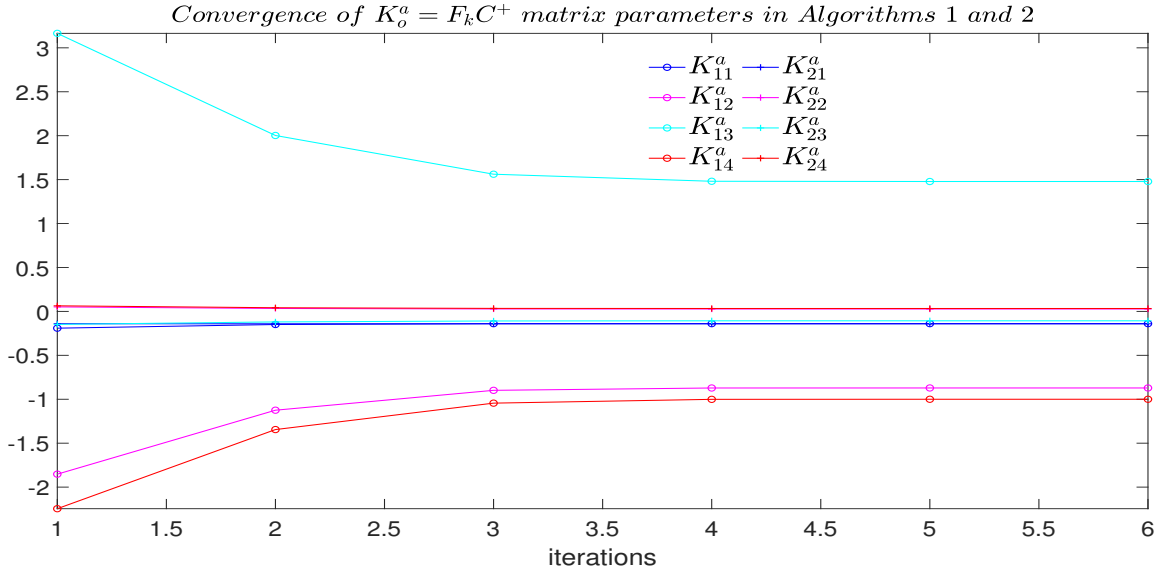


Figure 22: Convergence of the gain matrix  $K_o^a$  parameters by using the Algorithms 1 and 2.

Now, to examine the performance of the Algorithms 1 and 2, we set the initial sub-optimal stabilizing gain as

$$\mathbf{F}_0 = \begin{bmatrix} -0.0888 & -0.1875 & 0.7076 & -0.2328 \\ -0.1382 & 0.0105 & -0.0884 & 0.0141 \end{bmatrix} C. \quad (4.60)$$

The convergence of optimal gain matrix parameters  $K_o^a$  can be observed from Fig. 28. Note that their convergence properties are exactly the same as each other and they converge to the same gain matrix (4.58) as expected. Therefore, the output trajectories figure with this resultant optimal stabilizing gain  $K^a$  is the same as Fig. 27. Now compare the system responses in Fig. 26 and Fig. 27 to observe the regulation performance. The convergence of disturbance gain matrix parameters  $N_o$  is also illustrated in Fig. 30. Note that the converged parameters of  $N_o$  are the same the OPFB disturbance gain given in (4.59) as illustrated.

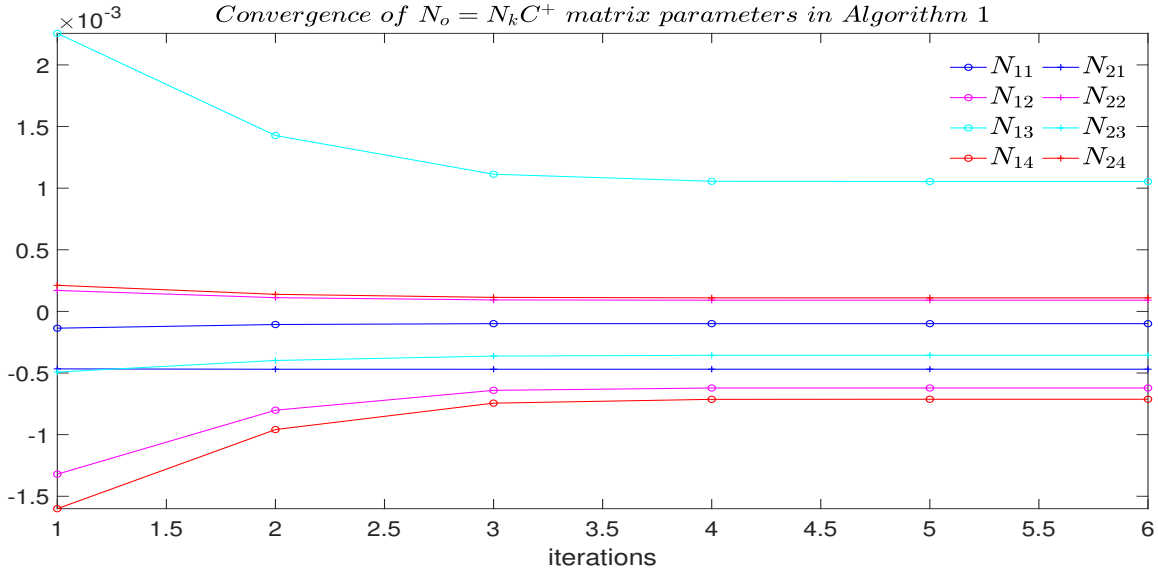


Figure 23: Convergence of the  $\mathbf{N}_o$  matrix parameters by using the Algorithm 1.

Next, to show the correct performance of model free Algorithm 3, we set the initial stabilizing gain  $\mathbf{F}_0$  as the same as (4.60). Then, based on the state information gathered from the system, the Algorithm 3 is employed to learn the Nash equilibrium gain solution  $\mathbf{K}_o^a$  online. Once the IRL Algorithm converged, we applied the corresponding control  $\mathbf{u} = -\mathbf{K}_o^a \mathbf{y}$  using to the actual system (4.1). The Algorithm 3 is also converged to the same Nash gain  $\mathbf{K}_o^a$  (4.58), thereby the resultant output vector states are the same as Fig. 27 as expected.

On the other hand, to check the correctness of the Algorithm 4, we select  $\mathbf{Q} = k\mathbf{C}^T\mathbf{C}$  with  $k = 0.05$  as explained in the Section 4.3.1. The resultant output trajectories are shown in the Fig. 6. Additionally, we have calculated the the Nash state feedback gain expression  $\mathbf{K}^*$  given in (5.20), and also the observability matrix  $\mathcal{O}$  for the system (4.1). Then, we compared the  $\mathbf{K}^*$  with the  $\bar{\mathbf{F}}_{k+1}\mathcal{O}$ . It has been seen that the two matrices have almost the same elements and  $L_2$  norm difference of them is calculated as 0.84, which verifies the correct performance of the Algorithm 4.

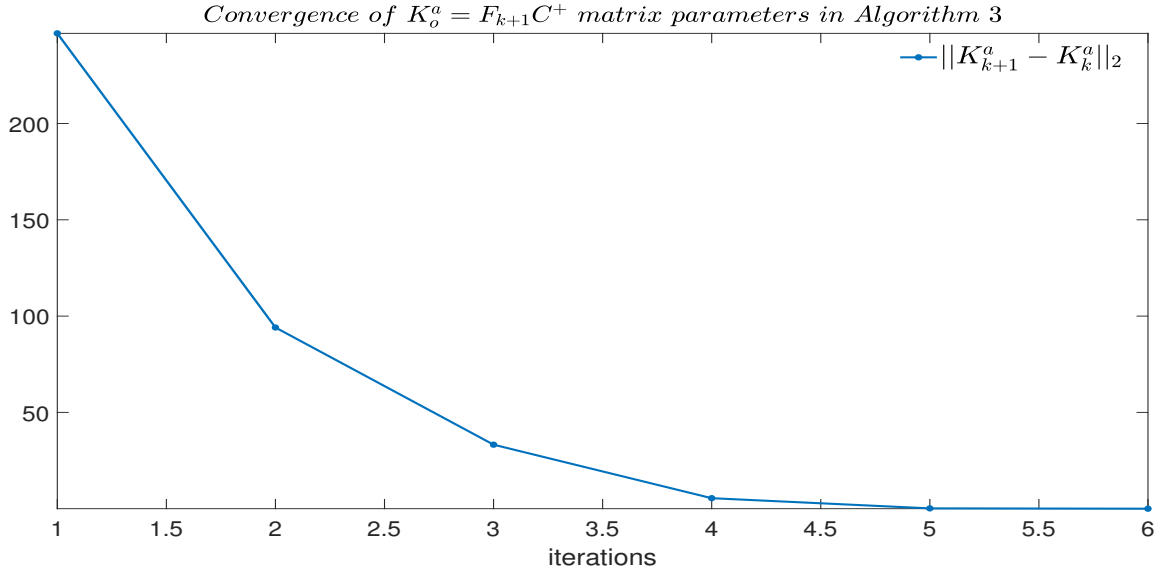


Figure 24: Convergence of the gain matrix  $\mathbf{K}_o^a$  parameters by using the Algorithm 3.

Lastly, the critical attenuation level obtained in the Theorems 6 and 11 is  $\gamma^* = 0.05$ . Note that with this critical attenuation  $\gamma^*$  level, the sufficient condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq 2(\gamma^*)^{-2}\mathbf{D}\mathbf{D}^T$  is relaxed. However, this condition indeed required for the Kleinman based Algorithms 1-4, and the critical attenuation level is obtained as  $\gamma^* = 0.49$ , which is also compatible with the sufficient condition  $\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \geq 2(\gamma^*)^{-2}\mathbf{D}\mathbf{D}^T$ . Therefore, we conclude that the model free algorithms reduces the  $L_2$  gain performance.

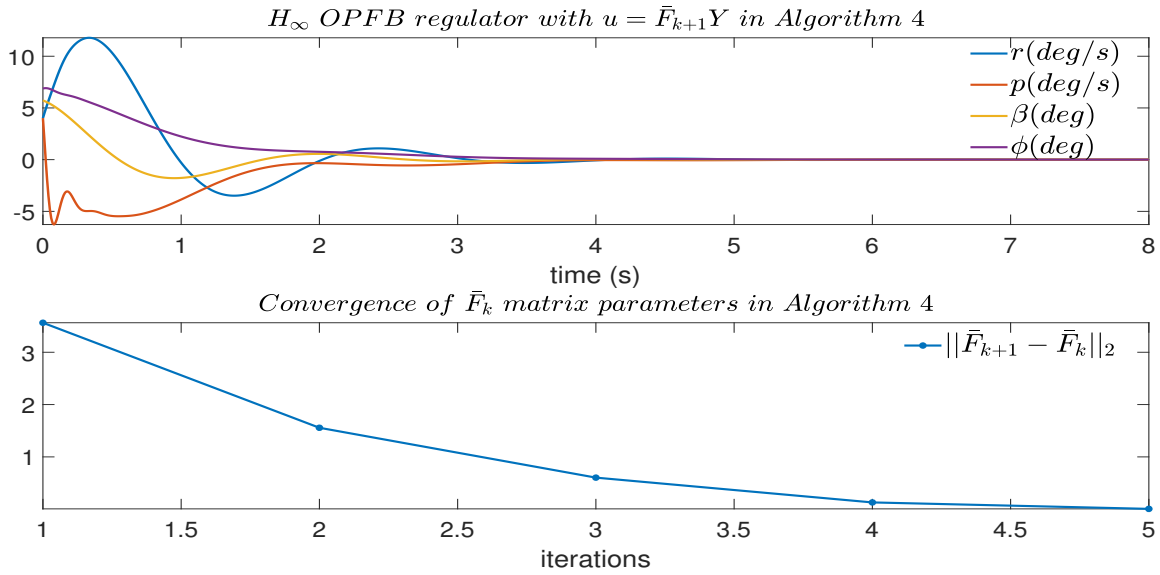


Figure 25: Performance of the Algorithm 4.

Adversarial Multi-agent Output Containment Graphical Game with Local and Global Objectives for UAVs.

Yusuf Kartal, Ahmet Taha Koru, Frank Lewis, Yan Wan, Atilla Dogan, 2021.

*Under Review in IEEE Transactions on Control of Network Systems.*

## CHAPTER 5

### Adversarial Multi-agent Output Containment Graphical Game with Local and Global Objectives for UAVs

Multi-agent systems (MAS) control is one of the most widely studied phenomena in recent years, due to the capability of MAS to perform certain tasks such as transportation [109], surveillance and reconnaissance [110], and target search and detection [111]. Three recognized categories of MAS control are single-leader & single-follower [112], single-leader & multi-follower [96], and multi-leader & multi-follower (MLF) [19]. For the MLF systems, one control objective is to guarantee the convergence of the output of each follower to the dynamic convex hull spanned by the outputs of leaders [19], [20], [21] which can be achieved by regulating the corresponding output containment error. These studies do not consider a practical case where the followers have mutual interests. Formulating the MLF system as a differential graphical game can address the mutual interests among the agents, and hence constitute a correct framework for the analysis of the MLF systems.

In the last decade, extensive efforts have been made to achieve Nash equilibrium strategies in differential graphical games [22],[23], [24], [25], [26]. The Nash equilibrium strategy in [22] is a centralized one since it requires to access the global state information of MAS. Min-max strategies in [76] guarantee a security-level performance when a Nash solution to graphical game does not exist. A modified objective functional in [23] guarantees the existence of both distributed and Nash solutions for the single-leader & multi-follower MAS.



In practice, besides Nash strategies, a MAS control design requirements may include  $\mathcal{L}_2$  gain boundedness, robustness against external disturbances, and output-feedback stabilizability. These design specifications are indeed local objectives of the graphical MLF game. Therefore, a new concept so called local and global objectives arises in the literature [113]. The works [76], [114], [115] are indeed related to the this new concept. [76] analyzed the multi-agent pursuit-evasion game where pursuer group have local objective, staying together, and global objective, capturing the evaders. [113] employed local objectives for a subset of agents, which on the other hand, are tasks determined around the global objective by agent-specific exosystems. Most of these works do not yield local and graphical Nash solutions simultaneously.

The standard existing solutions to the output containment MLF games uses an additive gain matrices to prove stability [19], [20]. In [19], sufficient local conditions in terms of stabilizing the local followers' dynamics and satisfying a certain  $\mathcal{H}_\infty$  criterion are investigated without considering mutual interests among the followers. The resulting strategies do not constitute a Nash equilibrium solution. In this paper, we propose graph theoretic and algebraic conditions to obtain Nash equilibrium strategies for the output containment of a MLF system. The salient contributions of this paper summed up into four categories as:

- The graphical output containment game is introduced where the mutual interests among the followers that experiences worst-case disturbances is addressed by introducing differential graphical games. For this graphical game, the output containment problem is the global objective of the follower group and disturbance attenuation is the local objective for some agents.
- A new solution to the  $\mathcal{H}_\infty$  static output feedback (OPFB) control problem that yields Nash equilibrium strategies for the graphical game is presented. The Nash solutions are proved to guarantee stability and bounded  $L_2$  gain considering

corresponding local min-max game worst-case disturbances with the developed novel necessary and sufficient conditions.

- The conditions for the distributed strategy based on available output data is presented to accomplish the local  $\mathcal{H}_\infty$  objective and the global output containment objective for the MLF system.
- To verify correctness of the proposed methods, the linear quadrotor model is developed around a particular flight condition. Then, the strategies constituting Nash equilibrium solution to output containment game are tested by means of the MLF quadrotor UAV simulations.

The rest of the paper is organised as follows. Section 5.2.1 considers a local  $\mathcal{H}_\infty$  game where the control inputs aim to minimize corresponding objective functional, whereas the disturbances aim to maximize. This corresponds to a well-known zero-sum game concept [61]. After playing the local game, we consider differential graphical game in Section 5.2.2 where followers' mutual interests are addressed. Furthermore, we give sufficient conditions on the objective functional's design parameters to satisfy both of the local and global objectives simultaneously. Section 5.3 reveals the  $\mathcal{L}_2$  gain bound stability analysis of the games introduced. Lastly, the proposed methods are tested by means of the MLF quadrotor UAV simulations in Section 5.4.

**Notations.** We use the following notations throughout this paper.  $\mathbf{I}_n \in \mathbb{R}^{n \times n}$  is the identity matrix,  $\mathbf{1}_N \in \mathbb{R}^N$  is a vector whose elements are all ones, and similarly  $\mathbf{0}_N \in \mathbb{R}^N$  is a vector whose elements are all zeros.  $\mathbf{0}_{n \times m} \in \mathbb{R}^{n \times m}$  is a matrix whose elements are all zeros.  $diag(\zeta_i)$  represents a diagonal matrix with  $\zeta_i \forall i \in 1, \dots, N$  on its diagonal. The condition  $A > 0$  ( $\geq 0$ ) denotes the positive (semi) definiteness of a matrix. The operator  $tr()$  denotes the trace of a matrix.  $\mathbf{C}^+ = \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}$  is the right-inverse of the full row-rank matrix  $\mathbf{C}$  and the Kronecker product operator is denoted by  $\otimes$ . The determinant of a square matrix is denoted by  $|\cdot|$ . Lastly, distance

from  $\mathbf{x} \in \mathbb{R}^n$  to the set  $\mathbb{C} \subseteq \mathbb{R}^n$  is defined via Euclidean norm as  $dist(\mathbf{x}, \mathbb{C}) = \inf_{\mathbf{y} \in \mathbb{C}} \|\mathbf{x} - \mathbf{y}\|_2$ .

## 5.1 Preliminaries

This section presents various definitions on adversarial multi-agent leader-follower (MLF) games by revealing multi-agent system dynamics and communication graph topologies that are of interest. Till the simulation section, the MLF game with generalized Linear Time Invariant (LTI) system dynamics are analyzed. Then, the specific LTI system dynamics are introduced for quadrotor UAVs in the Section 5.4.

### 5.1.1 Graph Topologies

A communication graph for  $N$  followers is denoted with a pair  $G^f = (V^f, E^f)$ , where the set  $V^f = \{v_1^f, \dots, v_N^f\}$  denotes the nodes, and  $E^f$  stands for the edges, which is composed of node pairs  $(v_i^f, v_k^f)$ . Each edge  $(v_k^f, v_i^f) \in E^f$ , has a weight  $a_{ik}^f = 1$  if node  $k$  is connected to node  $i$  and  $a_{ik}^f = 0$  otherwise. The graph is called as undirected if  $a_{ik}^f = a_{ki}^f, \forall i, k$ , otherwise it is termed a digraph. The in-degree matrix  $\mathbf{D}^f = diag\{d_i^f\}$  where  $d_i^f = \sum_{k=1}^N a_{ik}^f$ . The matrix  $\mathbf{A}^f = [a_{ik}^f] \in \mathbb{R}^{N \times N}$  denotes the adjacency or connectivity matrix. Then, the graph Laplacian matrix for the follower group is defined as  $\mathbf{L}^f = \mathbf{D}^f - \mathbf{A}^f$ . In this paper, we assume that  $G^f$  is a digraph that does not contain self loops, i.e,  $a_{ii}^f = 0$ .

In addition, define a bipartite graph  $G^b = (V^f, V^l, E^b)$  that consists of follower nodes  $V^f$ , leader nodes  $V^l$ , and edges  $E^b$ , which captures the information exchange among the follower and leader groups. Let  $g_{ji}^b$  denote the pinning of follower  $i$  to leader  $j$ , with  $g_{ji}^b = 1$  if follower  $i$  is connected to leader  $j$ , and  $g_{ji}^b = 0$  otherwise. Additionally, let  $g_{ij}^b$  denote the pinning of leader  $j$  to follower  $i$ , with  $g_{ij}^b = 1$  if leader  $j$  is pinned to follower  $i$ , and  $g_{ij}^b = 0$  otherwise. Then pinning matrices of follower  $i$

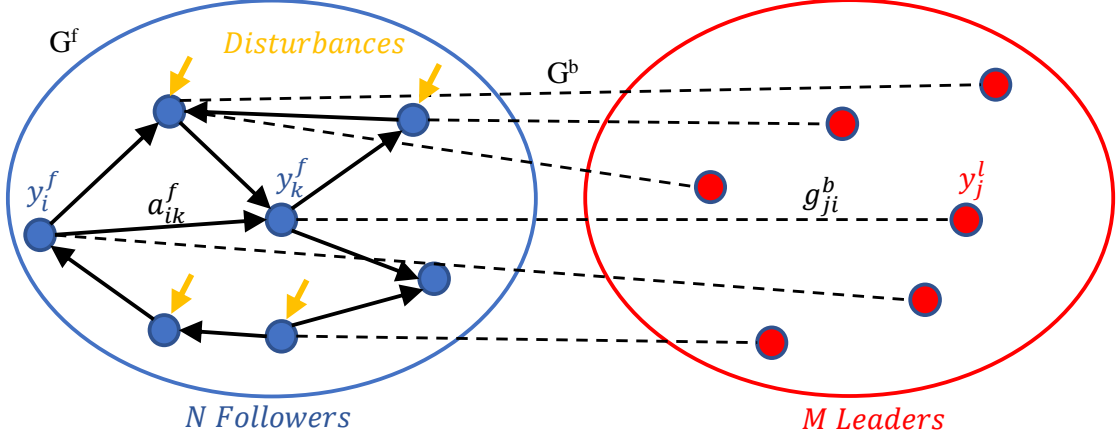


Figure 26: Layout of Adversarial Leader-Follower Graphical Game

and leader  $j$  in  $G^b$  are  $G_j^b$  and  $G_i^b$  respectively. One of the contributions of this paper is to analyze MLF games with directed communication graph topologies that lead to the limited measurement capabilities among the agents ( $G^f$  nodes) of game. Note that the graph formulation in this paper is valid for the non-negative connectivity and pinning weights i.e.,  $a_{ik}^f \geq 0$  and  $g_{ij}^f \geq 0$  but they are selected as ones and zeros so as not to increase gain of outer control loop for quadrotor UAVs in Section 5.4.

### 5.1.2 Multi-agent System Dynamics and Local Errors

Consider a follower group consisting of  $N$  agents with the dynamics given by

$$\begin{aligned} \dot{\mathbf{x}}_i^f &= \mathbf{A}\mathbf{x}_i^f + \mathbf{B}\mathbf{u}_i + \mathbf{D}\mathbf{w}_i \\ \mathbf{y}_i^f &= \mathbf{C}\mathbf{x}_i^f, \end{aligned} \tag{5.1}$$

$\forall i \in \{1, \dots, N\}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{D} \in \mathbb{R}^{n \times p}$  are system-state, input, disturbance matrices, and  $\mathbf{C} \in \mathbb{R}^{q \times n}$  is assumed to be a full row-rank output matrix to avoid redundant measurements. The corresponding vectors  $\mathbf{x}_i^f(t)$ ,  $\mathbf{u}_i(t)$ ,  $\mathbf{w}_i(t)$ , and  $\mathbf{y}_i^f(t)$  stand for the state, input, disturbance and output of  $i^{\text{th}}$  follower respectively.

Additionally, consider a leader group that consists of  $M$  agents with the dynamics given by

$$\begin{aligned}\dot{\mathbf{x}}_j^l &= \mathbf{A}\mathbf{x}_j^l \\ \mathbf{y}_j^l &= \mathbf{C}\mathbf{x}_j^l,\end{aligned}\tag{5.2}$$

$\forall j \in \{1, \dots, M\}$ , where vectors  $\mathbf{x}_j^l(t)$ , and  $\mathbf{y}_j^l(t)$  stand for the state, and output of  $j^{\text{th}}$  leader respectively.

**Assumption 6.** *The pair  $(\mathbf{A}, \mathbf{B})$  is stabilizable and the pair  $(\mathbf{A}, \mathbf{C})$  is detectable.*

**Assumption 7.** *The system (5.1) is OPFB stabilizable in the sense that the row-space of output matrix  $\mathbf{C}$  contains the sub-space that is spanned by the right eigenvectors corresponding to the unstable modes of  $\mathbf{A}$ .*

**Assumption 8.** *The non-zero columns of the output matrix  $\mathbf{C}$  are linearly independent.*

**Remark 16.** *The Assumption 7 can be interpreted such that all unstable modes are measured by the output matrix  $\mathbf{C}$  that represents the sensors installed in the systems (5.1) and (5.2). The Assumption 8 enables us to recover a state element  $x_i$  that is spanned in the row space of  $\mathbf{C}$  precisely from the output vector  $y$  once it is left multiplied with  $\mathbf{C}^+$ .*

**Definition 1.** *A set  $\mathbb{S}$  is convex if the line segment between any two points in  $\mathbb{S}$  lies in  $\mathbb{S}$ , i.e.,  $(\theta s_1 + (1 - \theta)s_2) \in \mathbb{S}$  holds for any  $s_1, s_2 \in \mathbb{S}$  and any  $\theta \in [0, 1]$  [116]. The convex hull of a set  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$  is denoted by  $\text{conv}(\mathcal{X}) = \{\sum_{j=1}^M \gamma_j \mathbf{x}_j | \mathbf{x}_j \in \mathcal{X}, \gamma_j \geq 0, \sum_{j=1}^M \gamma_j = 1\}$ .*

Let  $Y_{G^l} = \{\mathbf{y}_1^l, \dots, \mathbf{y}_M^l\}$  be a set of leaders' positions. Note that the centroid of leaders' positions  $\sum_{j=1}^M \gamma_j \mathbf{y}_j^l$  with  $\gamma_j = 1/M, \forall j \in \{1, \dots, M\}$  is an element of  $\text{conv}(Y_{G^l})$ .

**Definition 2.** Consider the  $i^{\text{th}}$  follower dynamics (5.1) and the  $j^{\text{th}}$  leader dynamics (5.2). The outputs of all followers are said to converge to the convex hull spanned by the outputs of leaders if

$$\lim_{t \rightarrow \infty} \text{dist}(\mathbf{y}_i^f, \text{conv}(Y_{Gl})) = 0, \forall i \in 1, \dots, N. \quad (5.3)$$

Motivated by [19], the local relevant output containment vector for  $i^{\text{th}}$  follower is defined as

$$\begin{aligned} \boldsymbol{\xi}_i &= \sum_{k=1}^N a_{ik}^f (\mathbf{y}_k^f - \mathbf{y}_i^f) + \sum_{j=1}^M g_{ij}^b (\mathbf{y}_j^f - \mathbf{y}_i^l) \\ &= \sum_{j=1}^M \left( \frac{1}{M} \sum_{k=1}^N a_{ik}^f (\mathbf{y}_k^f - \mathbf{y}_i^f) + g_{ij}^b (\mathbf{y}_j^f - \mathbf{y}_i^l) \right). \end{aligned} \quad (5.4)$$

The global form of (5.4) is

$$\boldsymbol{\xi} = - \sum_{j=1}^M \left( (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) (\mathbf{y}^f - \underline{\mathbf{y}}_j^l) \right) \quad (5.5)$$

where  $\boldsymbol{\xi} = [\boldsymbol{\xi}_1^T, \dots, \boldsymbol{\xi}_N^T]^T$ ,  $\boldsymbol{\psi}_j = \frac{1}{M} \mathcal{L}^f + \mathbf{G}_j^b$ ,  $\underline{\mathbf{y}}_j^l = \mathbf{1}_N \otimes \mathbf{y}_j^l$ , and  $\mathbf{y}^f = [\mathbf{y}_1^f, \dots, \mathbf{y}_N^f]^T$ .

The following steps enable us to verify that the output containment control objective (5.3) is achieved if  $\lim_{t \rightarrow \infty} \boldsymbol{\xi} = \mathbf{0}$ .

**Definition 3.** A matrix is called semi-simple if the geometric multiplicity of each eigenvalue is equal to its algebraic multiplicity.

It is known that a matrix is diagonalizable over real set  $\mathbb{R}$ , if it is semi-simple and every eigenvalue is real [117].

**Definition 4.** A follower is called a well-informed one if it is pinned to the all leaders and uninformed one if it is not pinned to any leader.

**Assumption 9.** The follower  $i$  in  $G^f$  is either a well-informed one or an uninformed one. There exists a directed path from one of the well informed followers  $k$  to an

uninformed follower  $i$  with  $k \neq i \forall k = \{1, \dots, N\}$  in  $G^f$ . Additionally, the graph is assumed to be acyclic and there is no directed connection from any uninformed follower  $i$  to well informed follower  $k$ .

Note that the existence of directed path from well-informed follower  $k$  to uninformed follower  $i$  is standard in multi-leader output containment control [19],[20]. Apart from this, we also assume that the graph  $G^f$  is acyclic and there is no directed connection from any uninformed follower  $i$  to well informed follower  $k$ , which allows us to modify the standard objective functionals in the MLF games. The next corollary is an essential step before we introduce further details on the MLF games.

**Remark 17.** Given Assumption 9,  $\psi_j$  and  $\sum_{j=1}^M \psi_j$  are non-singular  $M$  matrices, and hence the real parts of their eigenvalues are positive [118]. Furthermore, since the graph  $G^b$  is assumed to be acyclic, the pinned Laplacian matrix  $(\sum_{j=1}^M \psi_j)$  is indeed a simple matrix that has purely positive real eigenvalues, which is introduced in the Chapter 5.7.3 of [119].

**Remark 18.** Assumption 9 also implies that there exists a directed path from at least one of the leaders to any follower in a unified communication graph  $G^f \cup G^b$ . This allows us to verify non-singular  $M$ -matrix property of the matrices  $\psi_j$  and  $\sum_{j=1}^M \psi_j$ . Therefore, their inverses,  $(\psi_j)^{-1}$  and  $(\sum_{j=1}^M \psi_j)^{-1}$ , exist and are non-negative. A comprehensive fifty properties of  $M$ -matrices are detailed in [120].

**Lemma 3.** Given the Assumption 9, consider the follower and leader dynamics in (5.1) and (5.2). The condition (5.3) is achieved if  $\lim_{t \rightarrow \infty} \boldsymbol{\xi} = \mathbf{0}$ .

*Proof:* Re-write the global form of relevant output containment vector (5.5)

as

$$\boldsymbol{\xi} = - \left( \sum_{j=1}^M (\psi_j \otimes \mathbf{I}_q) \right) \mathbf{y}^f + \sum_{j=1}^M \left( (\psi_j \otimes \mathbf{I}_q) \underline{\mathbf{y}}_j^l \right) \quad (5.6)$$

$\lim_{t \rightarrow \infty} \boldsymbol{\xi} = \mathbf{0}$  implies that

$$\begin{aligned}
\mathbf{y}^f &\rightarrow \left( \sum_{j=1}^M (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) \right)^{-1} \sum_{j=1}^M \left( (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) \underline{\mathbf{y}}_j^l \right) \\
&\rightarrow \sum_{j=1}^M \left[ \left( \sum_{r=1}^M (\boldsymbol{\psi}_r \otimes \mathbf{I}_q) \right)^{-1} (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) \underline{\mathbf{y}}_j^l \right] \\
&\rightarrow \sum_{j=1}^M \left[ \left( \left( \sum_{r=1}^M \boldsymbol{\psi}_r \right)^{-1} \boldsymbol{\psi}_j \mathbf{1}_N \right) \otimes \underline{\mathbf{y}}_j^l \right]
\end{aligned} \tag{5.7}$$

as  $t \rightarrow \infty$ . Realize that

$$\begin{aligned}
\sum_{j=1}^M \left[ \left( \sum_{r=1}^M \boldsymbol{\psi}_r \right)^{-1} \boldsymbol{\psi}_j \mathbf{1}_N \right] &= \left( \sum_{r=1}^M \boldsymbol{\psi}_r \right)^{-1} \left( \sum_{j=1}^M \boldsymbol{\psi}_j \mathbf{1}_N \right) \\
&= \mathbf{1}_N,
\end{aligned} \tag{5.8}$$

which implies that each row sum of the vectors  $\left( \sum_{r=1}^M \boldsymbol{\psi}_r \right)^{-1} \boldsymbol{\psi}_j \mathbf{1}_N \forall j \in 1, \dots, M$  is equal to 1. Using Remarks 17 and 18, one can derive non-negative definiteness of each entry of the vectors  $\left( \sum_{r=1}^M \boldsymbol{\psi}_r \right)^{-1} \boldsymbol{\psi}_j \mathbf{1}_N$  [118]. Therefore,  $\lim_{t \rightarrow \infty} \boldsymbol{\xi} = \mathbf{0}$  indeed verifies that the objective (5.3) is achieved by Definition 1. This completes the proof.

□

### 5.1.3 Multi-agent Error Dynamics and $\mathcal{L}_2$ Gain Bound

Based on (5.7) in Lemma 3, define the global output containment vector as

$$\boldsymbol{\delta}_y = \mathbf{y}^f - \left( \sum_{j=1}^M (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) \right)^{-1} \sum_{j=1}^M \left( (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) \underline{\mathbf{y}}_j^l \right). \tag{5.9}$$



Note that  $\lim_{t \rightarrow \infty} \delta_{\mathbf{y}} = \mathbf{0} \implies \lim_{t \rightarrow \infty} \boldsymbol{\xi} = \mathbf{0}$  and (5.3) holds. Then, the global state containment vector is

$$\boldsymbol{\delta} = \mathbf{x}^f - \left( \sum_{j=1}^M (\boldsymbol{\psi}_j \otimes \mathbf{I}_n) \right)^{-1} \sum_{j=1}^M ((\boldsymbol{\psi}_j \otimes \mathbf{I}_n) \underline{\mathbf{x}}_j^l) \quad (5.10)$$

where  $\boldsymbol{\delta} = [\delta_1^T, \dots, \delta_N^T]^T$ ,  $\mathbf{x}^f = [\mathbf{x}_1^f, \dots, \mathbf{x}_N^f]^T$  and  $\underline{\mathbf{x}}_j^l = \mathbf{1}_N \otimes \mathbf{x}_j^l$ .

Consider the global containment vectors (5.9) and (5.10), and the follower dynamics in (5.1) and the leader dynamics in (5.2). Then, by using the fact that  $\boldsymbol{\xi} = - \left( \sum_{j=1}^M (\boldsymbol{\psi}_j \otimes \mathbf{I}_q) \right) \boldsymbol{\delta}_{\mathbf{y}}$ , the output containment error system can be written in the following global form

$$\begin{aligned} \dot{\boldsymbol{\delta}} &= (\mathbf{I}_N \otimes \mathbf{A})\boldsymbol{\delta} + (\mathbf{I}_N \otimes \mathbf{B})\mathbf{u} + (\mathbf{I}_N \otimes \mathbf{D})\mathbf{w} \\ \boldsymbol{\xi} &= - \left( \sum_{j=1}^M (\boldsymbol{\psi}_j \otimes \mathbf{C}) \right) \boldsymbol{\delta}. \end{aligned} \quad (5.11)$$

where  $\mathbf{u} = [\mathbf{u}_1^T, \dots, \mathbf{u}_N^T]^T$ , and  $\mathbf{w} = [\mathbf{w}_1^T, \dots, \mathbf{w}_N^T]^T$ .

**Definition 5.** Realization of the following inequality  $\forall \mathbf{w} \in [0, \infty)$  implies that the  $\mathcal{L}_2$  gain of system (5.11) is bounded by a prescribed disturbance attenuation level denoted by  $\gamma$

$$\int_0^\infty \frac{\boldsymbol{\xi}^T \boldsymbol{\xi}}{\|\mathbf{C}^T \mathbf{C}\|_2} dt \leq \gamma^2 \int_0^\infty \mathbf{w}^T \mathbf{w} dt + \beta \quad (5.12)$$

where  $\beta$  is a non-negative constant. The condition (5.12) is also called as nonexpansivity constraint in [74].

The following Lemma is inspired by [96].

**Lemma 4.** *Given the Assumption 9 and the output containment error system (5.11). Assume that a stabilizing static output feedback (OPFB) control takes the form*

$$\mathbf{u}_i = -c_i \mathbf{K}_i \hat{\boldsymbol{\xi}}_i. \quad (5.13)$$

where  $c$  is a positive constant. Then, the output containment objective (5.3) and  $\mathcal{L}_2$  gain bound condition (5.12) hold if the distributed  $N$  systems

$$\begin{aligned} \dot{\hat{\boldsymbol{\delta}}}_i &= \mathbf{A} \hat{\boldsymbol{\delta}}_i + \mathbf{B} \hat{\mathbf{u}}_i + \mathbf{D} \hat{\mathbf{w}}_i \\ &= (\mathbf{A} - c_i \lambda_i \mathbf{B} \mathbf{K}_i \mathbf{C}) \hat{\boldsymbol{\delta}}_i + \mathbf{D} \hat{\mathbf{w}}_i \\ \hat{\boldsymbol{\xi}}_i &= -\lambda_i \mathbf{C} \hat{\boldsymbol{\delta}}_i \end{aligned} \quad (5.14)$$

are both asymptotically stable and  $\mathcal{L}_2$  gain bounded by  $\gamma > 0$ , where  $\lambda_i$  are eigenvalues of the matrix  $\sum_{j=1}^M \boldsymbol{\psi}_j$ .

*Proof:* By Definition 3, the M-matrix  $\sum_{j=1}^M \boldsymbol{\psi}_j$  is diagonalizable given the Assumption 9, thereby there exists an invertible transformation matrix  $\mathbf{T}$  such that

$$\mathbf{T} \left( \sum_{j=1}^M \boldsymbol{\psi}_j \right) \mathbf{T}^{-1} = \boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N). \quad (5.15)$$

Let  $\hat{\boldsymbol{\delta}} = (\mathbf{T} \otimes \mathbf{I}_n) \boldsymbol{\delta}$ ,  $\mathbf{w} = (\mathbf{T}^{-1} \otimes \mathbf{I}_p) \hat{\mathbf{w}}$ , and  $\hat{\boldsymbol{\xi}} = (\mathbf{T} \otimes \mathbf{I}_q) \boldsymbol{\xi}$ . Then (5.14) can be written in the global form as

$$\begin{aligned} \dot{\hat{\boldsymbol{\delta}}} &= (\mathbf{I}_N \otimes \mathbf{A} - c_i \boldsymbol{\Lambda} \otimes \mathbf{B} \mathbf{K} \mathbf{C}) \hat{\boldsymbol{\delta}} + (\mathbf{I}_N \otimes \mathbf{D}) \hat{\mathbf{w}} \\ \hat{\boldsymbol{\xi}} &= -(\boldsymbol{\Lambda} \otimes \mathbf{C}) \hat{\boldsymbol{\delta}}. \end{aligned} \quad (5.16)$$

Realize that the  $\mathcal{L}_2$  norm equivalence of the pairs  $(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}})$  and  $(\mathbf{w}, \hat{\mathbf{w}})$  are straight forward as  $\mathbf{T}$  is invertible, and hence (5.11) and (5.16) have the same  $\mathcal{L}_2$  gains.

Additionally, since they have the same transfer function, stability of the equilibrium (origin) in both systems are equivalent, which completes the proof.  $\square$

## 5.2 Multi-agent leader-follower game formulation

In this section, two types of games are analyzed. The first game is a local game, which is played between the control and disturbance terms of the  $i^{th}$  follower. On the other hand, the second game is a global graphical game, which is played between the  $i^{th}$  and  $j^{th}$  followers considering the worst case disturbance that is achieved in the corresponding local game.

**Remark 19.** *Herein, we employ state feedback control method to derive the Nash equilibrium strategies. Then, the achieved control strategies can be re-written in the OPFB form given the Assumptions 6-8. This will be illustrated later in the Section 5.3.*

### 5.2.1 Local $\mathcal{H}_\infty$ Game Solution

This section presents a local  $\mathcal{H}_\infty$  game solution, which is formulated in the form of zero-sum game by introducing the following objective functional

$$\mathcal{J}_i(\hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i) = \int_0^\infty (\hat{\boldsymbol{\delta}}_i^T \mathbf{Q}_i \hat{\boldsymbol{\delta}}_i + \frac{1}{c_i \lambda_i} \hat{\mathbf{u}}_i^T \mathbf{R}_i \hat{\mathbf{u}}_i - \gamma^2 \hat{\mathbf{w}}_i^T \hat{\mathbf{w}}_i) d\tau \quad (5.17)$$

where  $\mathbf{Q}_i \geq 0$  and  $\mathbf{R}_i > 0$  are symmetric design matrices with appropriate dimensions. We assume that  $\mathbf{S}_i$  is selected such that the pair  $(\mathbf{A}, \sqrt{\mathbf{Q}_i})$  is observable.

Now the  $\mathcal{H}_\infty$  control problem can be solved by treating  $\hat{\mathbf{u}}_i$  as a minimizing player, whereas  $\hat{\mathbf{w}}_i$  as a maximizing player of the cost (5.17). Then, the game can be formulated as

$$\begin{aligned}\mathcal{V}_i(\hat{\boldsymbol{\delta}}_i) &\triangleq \mathcal{J}_i(\hat{\mathbf{u}}_i^*, \hat{\mathbf{w}}_i^*) = \min_{\hat{\mathbf{u}}_i} \max_{\hat{\mathbf{w}}_i} \mathcal{J}_i(\hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i) \\ &= \max_{\hat{\mathbf{w}}_i} \min_{\hat{\mathbf{u}}_i} \mathcal{J}_i(\hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i)\end{aligned}\quad (5.18)$$

where  $\mathcal{V}_i(\hat{\boldsymbol{\delta}}_i)$  denotes the value functional such that

$$\mathcal{V}_i(\hat{\boldsymbol{\delta}}_i) = \int_t^\infty (\hat{\boldsymbol{\delta}}_i^T \mathbf{Q}_i \hat{\boldsymbol{\delta}}_i + \frac{1}{c_i \lambda_i} \hat{\mathbf{u}}_i^T \mathbf{R}_i \hat{\mathbf{u}}_i - \gamma^2 \hat{\mathbf{w}}_i^T \hat{\mathbf{w}}_i) d\tau, \quad (5.19)$$

and the pair  $(\hat{\mathbf{u}}_i^*, \hat{\mathbf{w}}_i^*)$  is the local game theoretic saddle point (local game Nash equilibrium).

**Lemma 5.** *The pair  $(\hat{\mathbf{u}}_i^*, \hat{\mathbf{w}}_i^*)$  constitutes a Nash equilibrium of the game (5.18) if*

$$\hat{\mathbf{u}}_i^* = -c_i \lambda_i \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i \hat{\boldsymbol{\delta}}_i, \quad (5.20)$$

$$\hat{\mathbf{w}}_i^* = \gamma^{-2} \mathbf{D}^T \mathcal{P}_i \hat{\boldsymbol{\delta}}_i, \quad (5.21)$$

where  $\mathcal{P}_i$  is the solution of corresponding Game Algebraic Riccati Equation (GARE) such that

$$\mathcal{P}_i \mathbf{A} + \mathbf{A}^T \mathcal{P}_i - c_i \lambda_i \mathcal{P}_i \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i + \mathbf{Q}_i + \gamma^{-2} \mathcal{P}_i \mathbf{D} \mathbf{D}^T \mathcal{P}_i = \mathbf{0}. \quad (5.22)$$

*Proof:* Begin with deriving the Hamiltonian to solve the minimizing & maximizing extrema that satisfies the Nash condition (5.18) as

$$\begin{aligned}\mathcal{H}_i(\nabla \mathcal{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i) &= \hat{\boldsymbol{\delta}}_i^T \mathbf{Q}_i \hat{\boldsymbol{\delta}}_i + \frac{1}{c_i \lambda_i} \hat{\mathbf{u}}_i^T \mathbf{R}_i \hat{\mathbf{u}}_i - \gamma^2 \hat{\mathbf{w}}_i^T \hat{\mathbf{w}}_i \\ &+ \nabla \mathcal{V}_i^T (\mathbf{A} \hat{\boldsymbol{\delta}}_i + \mathbf{B} \hat{\mathbf{u}}_i + \mathbf{D} \hat{\mathbf{w}}_i) = 0\end{aligned}\quad (5.23)$$

where  $\nabla\mathcal{V}_i = \partial\mathcal{V}_i/\partial\hat{\delta}_i$  is the co-state vector and the boundary condition is  $\mathcal{V}_i(0) = 0$ . Using the quadratic form  $\mathcal{V}_i(\hat{\delta}_i) = \hat{\delta}_i^T \mathcal{P}_i \hat{\delta}_i$  where  $\mathcal{P}_i = \mathcal{P}_i^T$ , and applying the stationarity conditions  $\partial\mathcal{H}_i(\nabla\mathcal{V}_i^*, \hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i)/\partial\hat{\mathbf{u}}_i = \mathbf{0}$  and  $\partial\mathcal{H}_i(\nabla\mathcal{V}_i^*, \hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i)/\partial\hat{\mathbf{w}}_i = \mathbf{0}$  yields the optimal control and disturbance respectively as (5.20) and (5.21). Additionally, substituting (5.20) and (5.21) into (5.23) and equating resultant HJI equation to zero, gives the Riccati equation (5.22).

Note that the sign of Hessians,  $\partial^2\mathcal{H}_i(\nabla\mathcal{V}_i^*, \hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i)/\partial\hat{\mathbf{u}}_i^2 > 0$  and  $\partial^2\mathcal{H}_i(\nabla\mathcal{V}_i^*, \hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i)/\partial\hat{\mathbf{w}}_i^2 < 0$ , along with the unboundedness of limits  $\lim_{d \rightarrow \infty} \mathcal{J}_i(\hat{\mathbf{u}}_i^*, \hat{\mathbf{w}}_i)$ ,  $\lim_{\hat{u}_i \rightarrow \infty} \mathcal{J}_i(\hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i^*)$  indeed show that (5.20) and (5.21) are the global optimal minimizing and maximizing extrema respectively. This indeed verifies that the pair  $(\hat{\mathbf{u}}_i^*, \hat{\mathbf{w}}_i^*)$  denotes the Nash equilibrium, which completes the proof.  $\square$

**Remark 20.** *The main purpose of introducing the local game (5.18 is to find the worst case disturbance for the  $i^{\text{th}}$  follower (5.1). Since  $i^{\text{th}}$  player is responsible for acting according to not only the worst case disturbance but also strategy of the  $j^{\text{th}}$  player,  $\hat{\mathbf{u}}_i^*$  (5.20) is not the finalized control strategy yet.*

**Corrolary 2.** *Let the distributed  $N$  systems (5.14) experience the worst-case disturbance derived in (5.21), i.e.,  $\hat{\mathbf{w}}_i = \hat{\mathbf{w}}_i^*$ , and introduce a new matrix  $\mathbf{F} = \mathbf{A} + \gamma^{-2}\mathbf{D}\mathbf{D}^T\mathcal{P}_i$  to facilitate the analysis. Then, the transformed local output containment error dynamics become*

$$\begin{aligned}\dot{\hat{\delta}}_i &= (\mathbf{A} - c_i\lambda_i\mathbf{B}\mathbf{K}\mathbf{C} + \gamma^{-2}\mathbf{D}\mathbf{D}^T\mathcal{P}_i)\hat{\delta}_i \\ &= \mathbf{F}\hat{\delta}_i + \mathbf{B}\hat{\mathbf{u}}_i \\ \hat{\xi}_i &= -\lambda_i\mathbf{C}\hat{\delta}_i\end{aligned}\tag{5.24}$$

as following is the natural outcome of using transformations  $\hat{\delta} = (\mathbf{T} \otimes \mathbf{I}_n)\delta$  and  $\mathbf{w}^* = (\mathbf{T}^{-1} \otimes \mathbf{I}_p)\hat{\mathbf{w}}^*$  on (5.21) such that

$$\hat{\mathbf{w}}^* = \gamma^{-2}(\mathbf{I}_N \otimes \mathbf{D}^T \mathcal{P}_i)\hat{\delta}. \quad (5.25)$$

## 5.2.2 Global Graphical Game Solution

This section proposes a modified cost functional for the graphical game whose players are the nodes of  $G^f$ . The main challenge in graphical games, is to design distributed Nash equilibrium control strategies, which cannot be achieved with the traditional quadratic cost functional formulation [23]. Therefore, a modified cost functional that provides both Nash and distributed control strategies in the sense that each follower uses the state information of its own and neighbors can be defined such that

$$J_i(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) = \int_0^\infty \sum_{k=1}^N (\hat{\delta}_{ik}^T \mathbf{Q}_{ik} \hat{\delta}_{ik} + \frac{1}{c_i \lambda_i} \hat{\mathbf{u}}_i^T \mathbf{R}_i \hat{\mathbf{u}}_i - \frac{a_{ik}^f}{(c_i \lambda_i)^2} \hat{\mathbf{u}}_k^T \mathbf{R}_{ik} \hat{\mathbf{u}}_k) d\tau. \quad (5.26)$$

where  $\hat{\delta}_{ik} = [\hat{\delta}_i^T, \hat{\delta}_k^T]^T$ ,  $\mathbf{Q}_{ik} = [\tilde{\mathbf{Q}}_i, \mathbf{0}_{n \times n}; \mathbf{0}_{n \times n}, \hat{\mathbf{Q}}_{ik}]$ , and  $\hat{\mathbf{u}}_{-i}$  is the set of control strategies of the  $k^{th}$  player in  $G^f$  where  $k \neq i$ . Then, the corresponding game is defined as

$$V_i^g(\hat{\delta}_i, \hat{\delta}_{-i}) \triangleq J_i(\hat{\mathbf{u}}_i^g, \hat{\mathbf{u}}_{-i}^g) \leq J_i(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}^g) \quad (5.27)$$

where the N-tuple  $\{V_1^g(\hat{\delta}_1, \hat{\delta}_{-1}), \dots, V_N^g(\hat{\delta}_N, \hat{\delta}_{-N})\}$  is called the global Nash equilibrium outcome and the N-tuple  $\{\hat{\mathbf{u}}_1^g, \dots, \hat{\mathbf{u}}_N^g\}$  denotes the Nash equilibrium strategies of the game (5.27).

**Remark 21.** A distributed Nash equilibrium strategy can be achieved if the functional  $V_i^g(\hat{\delta}_i, \hat{\delta}_{-i})$  depends on the its argument  $\hat{\delta}_i$  solely, i.e.,  $V_i^g(\hat{\delta}_i, \hat{\delta}_{-i}) = V_i^g(\hat{\delta}_i)$ . Otherwise, the resultant game optimal strategy would not be a distributed strategy.

Define a quadratic function that will play a key role in the corresponding Hamiltonian derivation such that

$$V_i(\hat{\delta}_i, \hat{\delta}_{-i}) \triangleq V_i(\hat{\delta}_i). \quad (5.28)$$

By the Remark 21, the following holds

$$\begin{aligned} \frac{\partial \tilde{V}_i^T}{\partial \hat{\delta}_{ik}} \dot{\hat{\delta}}_{ik} &= \begin{bmatrix} \frac{\partial \tilde{V}_i^T}{\partial \hat{\delta}_i} & \frac{\partial \tilde{V}_i^T}{\partial \hat{\delta}_k} \end{bmatrix}^T \begin{bmatrix} \dot{\hat{\delta}}_i \\ \dot{\hat{\delta}}_k \end{bmatrix} \\ &= \frac{\partial \tilde{V}_i^T}{\partial \hat{\delta}_i} \left( \mathbf{F} \hat{\delta}_i + \mathbf{B} \hat{u}_i \right). \end{aligned} \quad (5.29)$$

**Theorem 9.** Given the cost functional (5.26), local output containment error dynamics and quadratic form (5.28). Assume that the element of matrix  $\mathbf{Q}_{ik}$  satisfies

$$\hat{\mathbf{Q}}_{ik} = a_{ik}^f \tilde{\mathbf{P}}_k \mathbf{B} \mathbf{R}_k^{-1} \mathbf{R}_{ik} \mathbf{R}_k^{-1} \mathbf{B}^T \tilde{\mathbf{P}}_k \quad (5.30)$$

where  $\tilde{\mathbf{P}}_i$  and  $\tilde{\mathbf{P}}_k$  solve the corresponding algebraic Riccati equations such that

$$\mathbf{F}^T \tilde{\mathbf{P}}_i + \tilde{\mathbf{P}}_i \mathbf{F} + \tilde{\mathbf{Q}}_i - c_i \lambda_i \tilde{\mathbf{P}}_i \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \tilde{\mathbf{P}}_i = \mathbf{0}, \quad (5.31)$$

$$\mathbf{F}^T \tilde{\mathbf{P}}_k + \tilde{\mathbf{P}}_k \mathbf{F} + \tilde{\mathbf{Q}}_k - c_k \lambda_k \tilde{\mathbf{P}}_k \mathbf{B} \mathbf{R}_k^{-1} \mathbf{B}^T \tilde{\mathbf{P}}_k = \mathbf{0}. \quad (5.32)$$

Then, the Nash equilibrium strategy for  $i^{\text{th}}$  follower in the graphical game (5.27) takes the following form

$$\hat{u}_i^g = -c_i \lambda_i \mathbf{R}_i^{-1} \mathbf{B}^T \tilde{\mathbf{P}}_i \hat{\delta}_i. \quad (5.33)$$

*Proof:* This proof consists of two parts. The first part deals with the selection of  $\mathbf{Q}_{ik}$  to satisfy (5.28) and the second part rigorously analyzes the Nash equilibrium property of minimizing controls  $\hat{\mathbf{u}}_i$  derived in the first part.

*Selecting  $\mathbf{Q}_{ik}$ :* By (5.29), the Hamiltonian associated with the cost functional (5.26) is

$$\begin{aligned} H_i(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) &= \sum_{k=1}^N \left( \nabla\tilde{V}_i^T (\mathbf{F}\hat{\boldsymbol{\delta}}_i + \mathbf{B}\hat{\mathbf{u}}_i) \right) \\ &+ \sum_{k=1}^N \left( \hat{\boldsymbol{\delta}}_{ik}^T \mathbf{Q}_{ik} \hat{\boldsymbol{\delta}}_{ik} + \frac{1}{c_i \lambda_i} \hat{\mathbf{u}}_i^T \mathbf{R}_i \hat{\mathbf{u}}_i - \frac{a_{ik}^f}{(c_i \lambda_i)^2} \hat{\mathbf{u}}_k^T \mathbf{R}_{ik} \hat{\mathbf{u}}_k \right) = 0 \end{aligned} \quad (5.34)$$

where  $\nabla\tilde{V}_i = \partial\tilde{V}_i/\partial\hat{\boldsymbol{\delta}}_i$  with the boundary condition  $\tilde{V}_i(\mathbf{0}) = \mathbf{0}$ . To find the best responses, check the stationarity condition  $\partial H_i(\nabla\tilde{V}_i^g, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i})/\partial\hat{\mathbf{u}}_i = \mathbf{0}$ , which yields

$$\hat{\mathbf{u}}_i^g = -\frac{c_i \lambda_i}{2} \mathbf{R}_i^{-1} \mathbf{B}^T \nabla\tilde{V}_i^g \quad (5.35)$$

Substituting (5.35) into (5.34) yields the Hamilton-Jacobi (HJ) equation, i.e.,  $H_i(\nabla\tilde{V}_i^g, \hat{\mathbf{u}}_i^g, \hat{\mathbf{u}}_{-i}^g) = 0$ . Additionally, using the quadratic form of  $\tilde{V}_i(\hat{\boldsymbol{\delta}}_i)$  (5.28) in the HJ equation, one obtains

$$\begin{aligned} H_i(\nabla\tilde{V}_i^g, \hat{\mathbf{u}}_i^g, \hat{\mathbf{u}}_{-i}^g) &= \\ &\sum_{k=1}^N \hat{\boldsymbol{\delta}}_i^T (\mathbf{F}^T \tilde{\mathbf{P}}_i + \tilde{\mathbf{P}}_i \mathbf{F} + \tilde{\mathbf{Q}}_i - c_i \lambda_i \tilde{\mathbf{P}}_i \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \tilde{\mathbf{P}}_i) \hat{\boldsymbol{\delta}}_k \\ &+ \sum_{k=1}^N \hat{\boldsymbol{\delta}}_k^T (\hat{\mathbf{Q}}_{ik} - a_{ik}^f \tilde{\mathbf{P}}_k \mathbf{B} \mathbf{R}_k^{-1} \mathbf{R}_{ik} \mathbf{R}_k^{-1} \mathbf{B}^T \tilde{\mathbf{P}}_k) \hat{\boldsymbol{\delta}}_k = 0. \end{aligned} \quad (5.36)$$

This gives the element of matrix  $\mathbf{Q}_{ik}$  (5.30), and the Riccati equations (5.31)-(5.32), which completes the first part of the proof.

*Graphical Nash Equilibrium:* In this part, we prove that the N-tuple  $\{\hat{\mathbf{u}}_1^g, \dots, \hat{\mathbf{u}}_N^g\}$  indeed constitutes the Nash equilibrium strategies of the game (5.27). Realize that substituting the quadratic form (5.28) into (5.35) yields (5.33).



Now re-write the Hamiltonian (5.34) as

$$H_i(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) = \sum_{k=1}^N H_{ik}(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) \quad (5.37)$$

where  $H_{ik}(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i})$  can be expressed by completing the squares as

$$\begin{aligned} H_{ik}(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) &= H_{ik}(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i^g, \hat{\mathbf{u}}_{-i}^g) \\ &+ \frac{1}{c_i\lambda_i} (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^g)^T \mathbf{R}_i (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^g) - \frac{a_{ik}^f}{(c\lambda_k)^2} (\hat{\mathbf{u}}_k - \hat{\mathbf{u}}_k^g)^T \mathbf{R}_{ik} (\hat{\mathbf{u}}_k - \hat{\mathbf{u}}_k^g) \\ &- \frac{2a_{ik}^f}{(c\lambda_k)^2} \hat{\mathbf{u}}_k^g{}^T \mathbf{R}_{ik} (\hat{\mathbf{u}}_k - \hat{\mathbf{u}}_k^g) \end{aligned} \quad (5.38)$$

since the equality  $\nabla\tilde{V}_i^T \mathbf{B}\mathbf{u}_i^g = -\frac{2}{c}(\hat{\mathbf{u}}_i^g)^T \mathbf{R}_i \hat{\mathbf{u}}_i^g$  holds by (5.35). Having the fact that the Hamiltonian (5.34) is a differential equivalent of the cost functional (5.26), we can write

$$J_i(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) = \int_{t_0}^{\infty} \sum_{k=1}^N H_{ik}(\nabla\tilde{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}) dt + V_i(\hat{\boldsymbol{\delta}}_i(t_0)). \quad (5.39)$$

According to rules of the game (5.27), select  $\hat{\mathbf{u}}_k = \hat{\mathbf{u}}_k^g$  and use the Hamiltonian form (5.38) in (5.39) to obtain

$$J_i(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}^g) = \int_{t_0}^{\infty} \sum_{k=1}^N \frac{1}{c_i\lambda_i} (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^g)^T \mathbf{R}_i (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^g) dt + V_i(\hat{\boldsymbol{\delta}}_i(t_0)), \quad (5.40)$$

which clearly satisfies the game condition  $J_i(\hat{\mathbf{u}}_i^g, \hat{\mathbf{u}}_{-i}^g) \leq J_i(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_{-i}^g)$ , and hence the N-tuple  $\{\hat{\mathbf{u}}_1^g, \dots, \hat{\mathbf{u}}_N^g\}$  is indeed the Nash equilibrium of the game (5.27). This completes the second part of the proof.  $\square$

**Remark 22.** Realize that the local  $\mathcal{H}_\infty$  Nash control strategy  $\hat{\mathbf{u}}_i^*$  (5.20) does not have the same form as the global graphical Nash control  $\hat{\mathbf{u}}_i^g$  (5.33). Therefore, we have not obtained a control strategy that belongs to Nash equilibrium for both games (5.18) and

(5.27). The next Theorem 10 gives the necessary conditions to provide a dual Nash control strategy by selecting appropriate  $\tilde{\mathbf{Q}}_i$  and  $\mathbf{R}_{ik}$  design matrices.

**Theorem 10.** The Nash equilibrium control strategy  $\hat{\mathbf{u}}_i^*$  (5.20) for the local  $\mathcal{H}_\infty$  game (5.18) stands for the same Nash equilibrium control strategy  $\hat{\mathbf{u}}_i^g$  (5.33) for the global graphical game (5.27) if the design matrices  $\tilde{\mathbf{Q}}_i$  and  $\mathbf{R}_{ik}$  are selected as

$$\tilde{\mathbf{Q}}_i = -\gamma^{-2}\mathcal{P}_i\mathbf{D}\mathbf{D}^T\mathcal{P}_i + \mathcal{Q}_i \quad (5.41)$$

$$\mathbf{R}_{ik} = N^2\mathbf{B}^T\mathcal{P}_i\tilde{\mathbf{Q}}_i\mathcal{P}_i\mathbf{B} + \tilde{\mathbf{R}}_{ik} \quad (5.42)$$

where  $\tilde{\mathbf{R}}_{ik} > 0$ .

*Proof:* First we work on selection of the design matrix  $\tilde{\mathbf{Q}}_i$ . Begin with expanding the Riccati equation (5.31) using  $\mathbf{F} = \mathbf{A} + \gamma^{-2}\mathbf{D}\mathbf{D}^T\mathcal{P}_i$ , which yields

$$\begin{aligned} \tilde{\mathbf{P}}_i\mathbf{A} + \mathbf{A}^T\tilde{\mathbf{P}}_i - c_i\lambda_i\tilde{\mathbf{P}}_i\mathbf{B}\mathbf{R}_i^{-1}\mathbf{B}^T\tilde{\mathbf{P}}_i + \tilde{\mathbf{Q}}_i \\ + \gamma^{-2}(\tilde{\mathbf{P}}_i\mathbf{D}\mathbf{D}^T\mathcal{P}_i + \mathcal{P}_i\mathbf{D}\mathbf{D}^T\tilde{\mathbf{P}}_i) = \mathbf{0}. \end{aligned} \quad (5.43)$$

Select  $\tilde{\mathbf{Q}}_i$  as in (5.41), and assume that the positive definite matrix  $\mathcal{P}_i$  solves the local  $\mathcal{H}_\infty$  GARE (5.22). Then, to verify that it also solves the graphical game Riccati equation (5.43), it is required to obtain the equality  $\tilde{\mathbf{P}}_i = \mathcal{P}_i$  where  $\tilde{\mathbf{P}}_i$  is the solution of Graphical Riccati equation (5.43), illustrated in the following realization

$$\begin{aligned} (\tilde{\mathbf{P}}_i - \mathcal{P}_i)\mathbf{D}\mathbf{D}^T(\tilde{\mathbf{P}}_i - \mathcal{P}_i) = \tilde{\mathbf{P}}_i\mathbf{D}\mathbf{D}^T\tilde{\mathbf{P}}_i + \mathcal{P}_i\mathbf{D}\mathbf{D}^T\mathcal{P}_i \\ - \tilde{\mathbf{P}}_i\mathbf{D}\mathbf{D}^T\mathcal{P}_i - \mathcal{P}_i\mathbf{D}\mathbf{D}^T\tilde{\mathbf{P}}_i. \end{aligned} \quad (5.44)$$

Thereby (5.43) can be expressed in terms of  $\tilde{\mathbf{P}}_i$  with selected  $\tilde{\mathbf{Q}}_i$  (5.41) as

$$\tilde{\mathbf{P}}_i\mathbf{A} + \mathbf{A}^T\tilde{\mathbf{P}}_i - c_i\lambda_i\tilde{\mathbf{P}}_i\mathbf{B}\mathbf{R}_i^{-1}\mathbf{B}^T\tilde{\mathbf{P}}_i + \mathcal{Q}_i + \gamma^{-2}\tilde{\mathbf{P}}_i\mathbf{D}\mathbf{D}^T\tilde{\mathbf{P}}_i = \mathbf{0}. \quad (5.45)$$

Realize that (5.45) has a dual form of (5.22). Therefore, the Nash equilibrium control strategies  $\hat{\mathbf{u}}_i^*$  (5.20), and  $\hat{\mathbf{u}}_i^g$  (5.33) are the same as each other since the Riccati equation (5.22) has a unique positive definite solution under some conditions that will be detailed in Section 5.3. Additionally, setting the disturbance matrix  $\mathbf{D} = \mathbf{0}$  guarantees the equality  $\tilde{\mathbf{P}}_i = \mathbf{P}_i$  as  $\tilde{\mathbf{Q}}_i$  becomes  $\mathbf{Q}_i$  and (5.43) reduces to (5.45).

To complete the proof, the weighting matrix  $\mathbf{Q}_{ik}$  in (5.26) should be at least positive semi-definite with the selected design matrices (5.30) and (5.41)-(5.42). This can be achieved if both  $\tilde{\mathbf{Q}}_i$  and  $\hat{\mathbf{Q}}_{ik}$  are at least positive semi-definite matrices by applying the Schur complement method. Since  $\hat{\mathbf{Q}}_{ik}$  is clearly positive semi-definite, we conclude that  $\mathbf{Q}_i$  must be selected large enough to make  $\tilde{\mathbf{Q}}_i$  positive semi-definite as well. This completes the proof.  $\square$

**Remark 23.** *Theorem 10 implies that  $i^{\text{th}}$  player can minimize its global cost (5.26) by only playing the local  $\mathcal{H}_\infty$  (5.18) game. Hence, if the  $i^{\text{th}}$  player stabilizes the local error dynamics (5.14) with the Nash equilibrium control strategy (5.20), then both of the local and global objectives can be satisfied simultaneously.*

### 5.3 Stability and $\mathcal{L}_2$ gain bound analysis with output feedback

This section re-formulates the control strategy (local and global) derived in the Section 5.2 using the static output feedback method, and establishes corresponding  $\mathcal{L}_2$  gain bound by a prescribed attenuation level.

**Remark 24.** *Given the design parameters (5.30) and (5.41)-(5.42), the local  $\mathcal{H}_\infty$  game (5.18) solution and global graphical game (5.27) are indeed the same as each other by Theorem 10. Therefore, the stability and  $\mathcal{L}_2$  gain bound analysis conducted in this section, is valid for both games.*

Now, realize that the OPFB of  $i^{\text{th}}$  follower's control (5.13) can be re-written in the transformed coordinates as

$$\hat{\mathbf{u}}_i = -c_i \lambda_i \mathbf{K}_i \mathbf{C} \hat{\boldsymbol{\delta}}_i. \quad (5.46)$$

The next main theorem is a key to find the OPFB matrix  $\mathbf{K}_i^*$ .

**Theorem 11.** *Given Assumptions 6-8, the undisturbed system (5.14), i.e.,  $\hat{\mathbf{w}}_i = \mathbf{0}$ , is asymptotically stable if*

$$c_i = \frac{1}{\lambda_i}; \quad \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \geq \gamma^{-2} \mathbf{D} \mathbf{D}^T \quad (5.47)$$

where  $\gamma^*$  is the critical attenuation level [102]. Furthermore, the disturbed system (5.14) is  $L_2$  gain bounded by  $\gamma$  and asymptotically stable using the control  $\hat{\mathbf{u}}_i = c \mathbf{K}_i^* \hat{\boldsymbol{\xi}}_i$  (5.46) and the disturbance  $\hat{\mathbf{w}}_i = -\mathbf{N}_i^* \hat{\boldsymbol{\xi}}_i$  where

$$\mathbf{K}_i^* = \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i \mathbf{C}^+ \quad (5.48)$$

$$\mathbf{N}_i^* = \gamma^{-2} \mathbf{D}^T \mathcal{P}_i \mathbf{C}^+ \quad (5.49)$$

if the following sufficient condition is satisfied

$$c_i = \frac{1}{\lambda_i}; \quad \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \geq 2\gamma^{-2} \mathbf{D} \mathbf{D}^T. \quad (5.50)$$

*Proof:* Herein, we first derive full state feedback gains and then explain how to apply them to the OPFB design. The proof consists of two parts, we first prove the  $L_2$  gain bound condition (5.12) and then asymptotic stability of the equilibrium origin considering the system (5.14) without and with the worst-case disturbance.

*L<sub>2</sub> gain bound condition:* Select  $c_i$  as given in (5.50) and re-write the Hamiltonian (5.23) by completing the squares as

$$\begin{aligned} \mathcal{H}_i(\nabla\mathcal{V}_i, \hat{\mathbf{u}}_i, \hat{\mathbf{d}}_i) &= \mathcal{H}_i(\nabla\mathcal{V}_i, \hat{\mathbf{u}}_i^*, \hat{\mathbf{d}}_i^*) + (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^*)^T R_i (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^*) \\ &\quad - \gamma^{-2} (\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*)^T (\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*). \end{aligned} \quad (5.51)$$

Note that the HJI equation  $\mathcal{H}_i(\nabla\mathcal{V}_i, \hat{\mathbf{u}}_i^*, \hat{\mathbf{d}}_i^*) = 0$  holds with the boundary condition  $\mathcal{V}_i(\mathbf{0}) = 0$ . Then, the objective functional (5.17) can be re-expressed as

$$\begin{aligned} \mathcal{J}_i(\hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i) &= \int_0^\infty \left( (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^*)^T R_i (\hat{\mathbf{u}}_i - \hat{\mathbf{u}}_i^*) \right. \\ &\quad \left. - \gamma^{-2} (\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*)^T (\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*) \right) dt + \mathcal{V}_i(\hat{\boldsymbol{\delta}}_i(0)). \end{aligned} \quad (5.52)$$

Realize that the condition  $\mathcal{J}_i(\hat{\mathbf{u}}_i, \hat{\mathbf{w}}_i) \leq \beta$  implies that the non-expansivity constraint (5.12) holds. Select  $\beta = \mathcal{V}_i(\hat{\boldsymbol{\delta}}_i(0))$ ,  $\hat{\mathbf{u}}_i = \hat{\mathbf{u}}_i^*$ . Then, (5.52) reduces to

$$\begin{aligned} \mathcal{J}_i(\hat{\mathbf{u}}_i^*, \hat{\mathbf{w}}_i) &= - \int_0^\infty \gamma^{-2} (\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*)^T (\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*) dt + \mathcal{V}_i(\hat{\boldsymbol{\delta}}_i(0)) \\ &\leq \mathcal{V}_i(\hat{\boldsymbol{\delta}}_i(0)), \quad \forall \hat{\mathbf{d}}_i \in (0, \infty). \end{aligned} \quad (5.53)$$

This proves that the  $L_2$  gain bound condition (5.12) holds with  $\beta = \mathcal{V}_i(\hat{\boldsymbol{\delta}}_i(0))$ ,  $\hat{\mathbf{u}}_i = \hat{\mathbf{u}}_i^*$ . Additionally, the value of game (5.18) with  $\hat{\mathbf{u}}_i = \hat{\mathbf{u}}_i^*$  and  $\hat{\mathbf{d}}_i = \hat{\mathbf{d}}_i^*$  is  $\mathcal{V}_i(\hat{\boldsymbol{\delta}}_i(0))$  by (5.52).

*Asymptotic stability (5.14):* To verify asymptotic stability of the undisturbed system (5.14), we consider the algebraic Riccati equation (5.45) and benefit from gain margin  $[c_{lower}, \infty)$  with  $c_{lower} < \frac{1}{2}$  property of the  $\mathcal{H}_\infty$  control[74]. Note that the  $\mathcal{H}_\infty$  has gain margin less than  $\frac{1}{2}$  by Chapter 10 in [74] but the lower bound  $c_{lower}$  is not precisely defined, and hence the conditions given in (5.47)-(5.50) are only sufficient.

With this in mind, if  $\mathbf{K}_i^* \mathbf{C} = \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i$  and  $c_i = \frac{1}{\lambda_i}$ , the control (5.46) becomes  $\hat{\mathbf{u}}_i = -\mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i \hat{\boldsymbol{\delta}}_i$ .

Now, consider the Lyapunov function candidate  $\mathcal{V}(\hat{\boldsymbol{\delta}}_i) = \hat{\boldsymbol{\delta}}_i^T \mathcal{P}_i \hat{\boldsymbol{\delta}}_i$ . Realize that  $\mathcal{P}_i$  is the unique positive definite solution of GARE (5.22) if the condition  $\mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \geq \gamma^{-2} \mathbf{D} \mathbf{D}^T$  is satisfied. This is illustrated in the following realization

$$\begin{aligned} \dot{\mathcal{V}} &= \dot{\hat{\boldsymbol{\delta}}_i}^T \mathcal{P}_i \hat{\boldsymbol{\delta}}_i + \hat{\boldsymbol{\delta}}_i^T \mathcal{P}_i \dot{\hat{\boldsymbol{\delta}}_i} \\ &= \hat{\boldsymbol{\delta}}_i^T \left( -\mathcal{P}_i \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i - \mathcal{Q}_i + \gamma^{-2} \mathcal{P}_i \mathbf{D} \mathbf{D}^T \mathcal{P}_i \right) \hat{\boldsymbol{\delta}}_i \\ &\leq -\hat{\boldsymbol{\delta}}_i^T \mathcal{Q}_i \hat{\boldsymbol{\delta}}_i \quad \Leftarrow \quad \mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \geq \gamma^{-2} \mathbf{D} \mathbf{D}^T. \end{aligned} \quad (5.54)$$

Then, the observability of  $(\mathbf{A}, \sqrt{\mathcal{Q}_i})$  verifies that the undisturbed system (5.14) is asymptotically stable. Notice that the condition  $\mathbf{B} \mathbf{R}_i^{-1} \mathbf{B}^T \geq \gamma^{-2} \mathbf{D} \mathbf{D}^T$  only enables designer to prove asymptotic stability of the closed-loop matrix  $(\mathbf{A} - c_i \lambda_i \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i)$  by solving GARE (5.22) [102]. To prove stability of the closed-loop matrix  $(\mathbf{A} - c_i \lambda_i \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i + \gamma^{-2} \mathbf{D}^T \mathcal{P}_i \hat{\boldsymbol{\delta}}_i)$  that considers  $\hat{\mathbf{w}}_i = \hat{\mathbf{w}}_i^*$  (5.21), use the gain margin approach. Herein the disturbance gain  $\mathbf{N}_i^* \mathbf{C} = \gamma^{-2} \mathbf{D}^T \mathcal{P}_i$  is employed. According to the gain margin approach, the sufficient condition (5.50) is immediate. Furthermore, since the upper bound for gain margin is infinite, the strict equality conditions in (5.47)-(5.50) can be relaxed as inequality conditions such that  $c_i \geq \frac{1}{\lambda_i}$ .

Lastly, we explain how to apply state feedback gains to the OPFB design. Note that if  $\mathbf{C}$  is an invertible matrix, then all local errors (5.14) would be regulated optimally as  $\mathbf{C}^+$  becomes  $\mathbf{C}^{-1}$ , and  $\mathbf{K}_i^* \mathbf{C} = \mathbf{R}_i^{-1} \mathbf{B}^T \mathcal{P}_i \mathbf{C}^{-1} \mathbf{C}$  with  $\mathbf{K}_i^*$  in (5.48). On the other hand, if it is not square but full row-rank, then only the states spanned by row space of the output matrix  $\mathbf{C}$  would be regulated optimally given Assumptions 7-8, and the fact that  $\mathbf{C}^+ \mathbf{C}$  projects  $\mathbb{R}^n$  onto the row space of  $\mathbf{C}$ . Additionally, the other states would converge to the origin given Assumption 7. Realize that Assumption 6

implies that there could be an unstable mode that is observable but does not belong to the row space of  $\mathbf{C}$ . Therefore, Assumptions 7 and 8 are indeed required to apply static state feedback gains to the OPFB design. Lastly, the system (5.14) is stable against the worst-case disturbance  $\hat{\mathbf{w}}_i = -\mathbf{N}_i^* \hat{\boldsymbol{\xi}}_i$  with  $\mathbf{N}_i^*$  in (5.49), which affects only the states that belong to the row space of  $\mathbf{C}$  given the condition (5.50). This completes the proof.  $\square$

**Remark 25.** *Note that the condition  $\mathbf{B}\mathbf{R}_i^{-1}\mathbf{B}^T \geq 2\gamma^{-2}\mathbf{D}\mathbf{D}^T$  is only sufficient condition for the stability of the disturbed system i.e., there may be a Nash gain solution which stabilizes (5.14) but does not satisfy  $\mathbf{B}\mathbf{R}_i^{-1}\mathbf{B}^T \geq 2\gamma^{-2}\mathbf{D}\mathbf{D}^T$ . This is because of the critical attenuation level  $\gamma^*$  that is introduced in [102]. According to the Chapter 2 of [102], the undisturbed system (5.14) is asymptotically stable when  $\gamma > \gamma^*$ . However, in that case, one may not achieve a positive definite solution  $P$  for the GARE (5.22), which should be avoided by Theorem 10. On the other hand, Assumptions 6-9 are indeed necessary for the MLF game to have a stabilizing Nash solution.*

**Remark 26.** *Note that the proposed methods in this paper aim to regulate the output containment error vector  $\boldsymbol{\xi}$  (5.5) instead of global state vector  $\boldsymbol{\delta}$  (5.11). The reasoning behind this approach is not to deal with inner loop of the system dynamics of the agents for the graphical game (5.27) since they can be treated separately from the outer loop. An example for the inner loop is the attitude control loop of the Unmanned Aerial Vehicles (UAVs) [5]. Interested reader can examine our work [27] to observe how to design a controller to deal with the attitude and position control loops of the multi-quadrotor systems separately. Herein, we only aim to control the position of the UAVs that correspond to the outer position loop, and the attitude controller is assumed be stable by Assumption 7.*

**Remark 27.** *If the graph  $G^f$  is a tree with the root node being a well informed follower  $k$  and each follower's connectivity weight in  $\mathcal{A}^f$  is multiplied with  $M$  and*

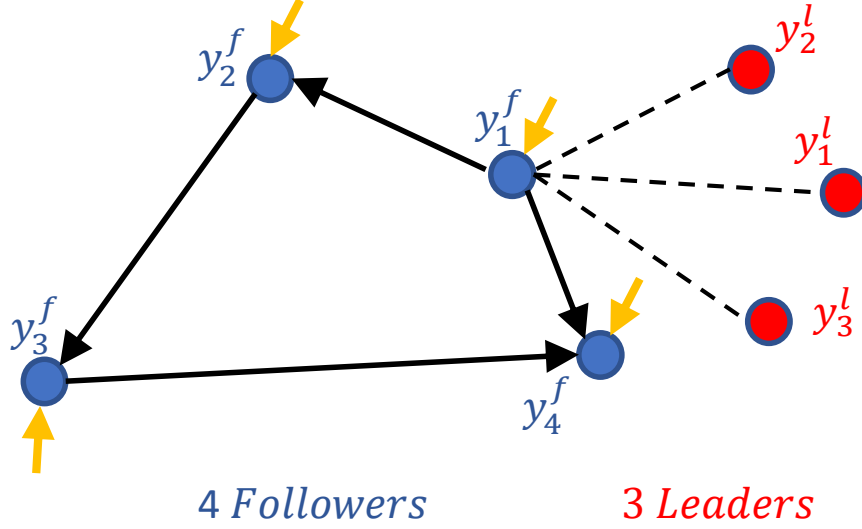


Figure 27: The communication graph of the four followers and three leaders MLF game.

all pinning gains are multiplied with  $1/M$ , the pinned Laplacian matrix  $(\sum_{j=1}^M \psi_j)$  would have all eigenvalues ones since the Laplacian matrix of the graph  $G^f$ , i.e.,  $\mathcal{L}^f$  is nilpotent matrix [121] in this case. Then, the term  $\lambda_i$  can be removed from the objective functionals (5.17) and (5.26), and hence from the corresponding Nash strategies (5.20) and (5.33). Therefore, the condition (5.47) would be modified as  $c = 1$ .

#### 5.4 Simulation results

In this section, we first introduce the multi-agent quadrotor Unmanned Aerial Vehicle (UAV) dynamics by using the backstepping control phenomena from our previous works [5] and [27]. Then, we apply the proposed methods to the MLF UAV game to verify correct performances of them.



It is well-known that the quadrotor linear model can be obtained around a flight condition that is called as hover where quadrotor's Body frame is aligned with the Earth (Inertial) frame (See Fig. 1 in [27] for frame illustrations and North-West-Up coordinate axis system convention). Then, the system matrices for linear model of the quadrotor's outer position & velocity control loop can be obtained as

$$\dot{\mathbf{x}}_i^f = \begin{bmatrix} \mathbf{0}_{3 \times 3} & \mathbf{I}_3 \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \end{bmatrix} \mathbf{x}_i^f + \begin{bmatrix} 0_3 & 0_3 & 0_3 & 0_3 \\ g & 0 & 0 & 0 \\ 0 & -g & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{m} \end{bmatrix} \mathbf{u}_i^f \quad (5.55)$$

where  $\mathbf{x}_i^f = [x_i \ y_i \ z_i \ u_i \ v_i \ w_i]^T$  is a vector of stacked position and velocity vectors,  $\mathbf{u}_i^f = [\theta_i \ \phi_i \ \psi_i \ \mu_i]^T$  is a vector of stacked Euler angles and total thrust produced by the quadrotor,  $g = 9.81m/sec^2$  is the gravitational acceleration and  $m = 0.467kg$  is the mass of the quadrotor UAV. To test the robustness, the disturbance matrix is selected as  $\mathbf{D} = [0 \ 0 \ 0 \ 0 \ 0 \ 1]^T$ . In this case, the output matrix is an identity matrix, which means that the follower group is required to measure only position and velocity information from the leader group, and hence there is no need for the attitude information of any agent in the MLF UAV game. Notice that Assumptions 7 and 8 are naturally satisfied with this well partitioned UAV model.

The nonstandard backstepping control method in [27] aims to find the desired Euler angles that needs to be tracked by the attitude controller. Herein, we use the same convention but assume that the  $\theta$  pitch angle and  $\phi$  roll angle remains in the interval  $[-\pi/6, \pi/6]$  to make sure that the quadrotor UAV shows a linear behaviour as small angle approximation is valid in this interval [5]. Furthermore, since the linear model of quadrotor is obtained around the hover flight condition, the input

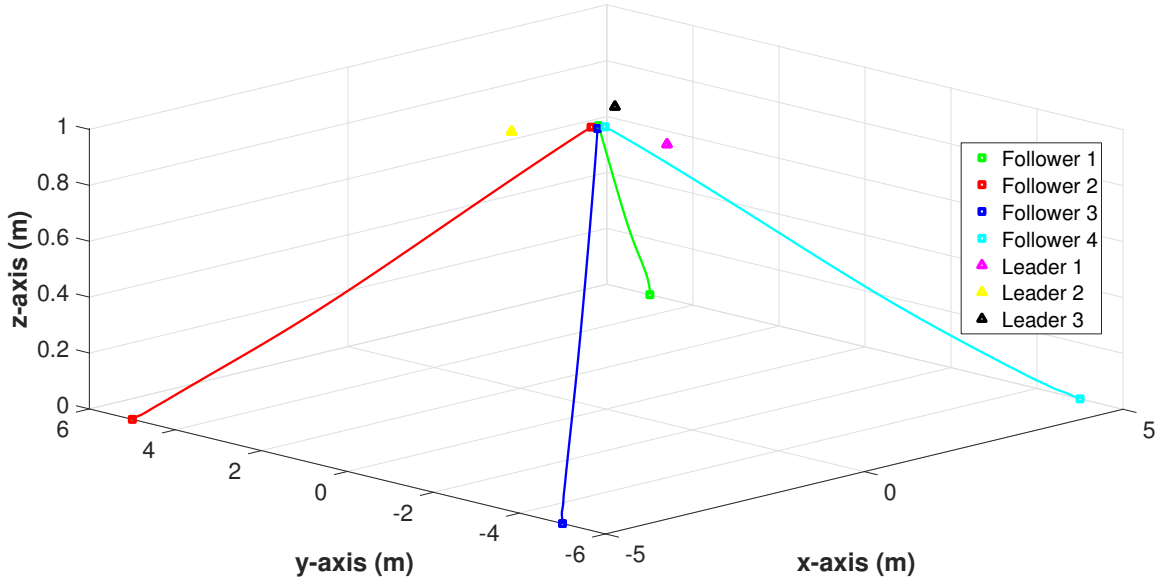


Figure 28: Path tracked by the UAVs in the MLF game.

parameters of (5.55) should be summed with the nominal flight condition parameters, i.e.,  $\mathbf{u}_{hover} = [\theta = 0 \ \phi = 0 \ \psi = c_\psi \ \mu = mg]^T$  before using them as set points for the attitude control loop. The attitude controller is taken from our work [27], and the proof of stability for the backstepping method is omitted.

The adjacency matrix  $\mathcal{A}^f$  that connects  $N = 4$  followers or nodes of the graph  $G^f$  and the bipartite communication graph parameters  $G_j^b \ j \in \{1, \dots, M = 3\}$  that do not contradict with Assumption 9 are selected as

$$\mathcal{A}^f = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}, \quad g_{11}^b = g_{12}^b = g_{13}^b = 1. \quad (5.56)$$

The MLF game communication graph is illustrated in Fig. 27. Note that the root node is selected as the first follower  $\mathbf{y}_1^f$  according to the Assumption 9. There

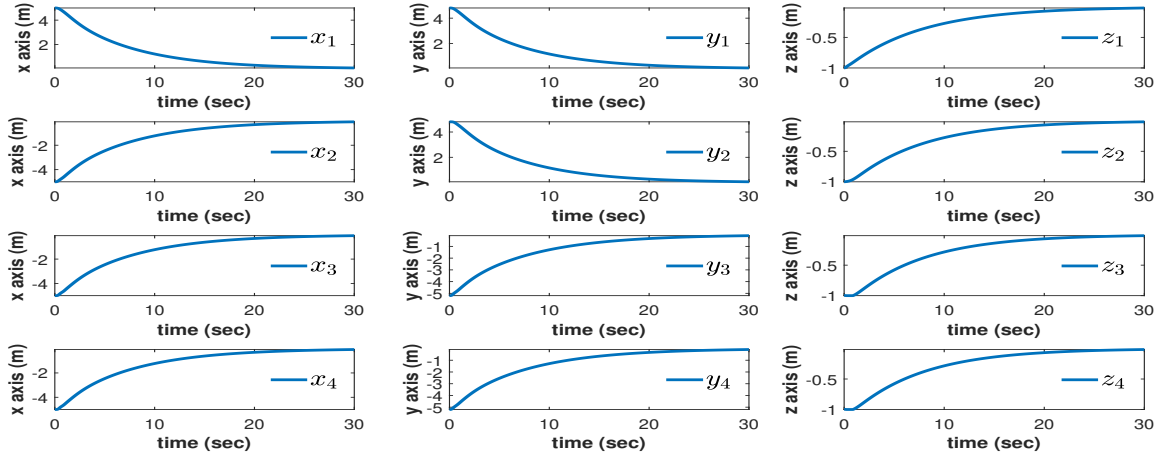


Figure 29: Output containment position error for each follower in three dimensional space.

is also no connection from uninformed followers ( $\mathbf{y}_2^f, \mathbf{y}_3^f, \mathbf{y}_4^f$ ) to the well informed follower  $\mathbf{y}_1^f$  and the follower graph is an acyclic digraph. The eigenvalues of the pinned Laplacian matrix ( $\sum_{j=1}^M \psi_j$ ) are 3.000, 0.667, 0.333, 0.333 that are all positive real numbers as required.

Now, to verify the correct performance of the proposed methods, the control for each follower is selected as (5.13) and the worst-case disturbance is set to  $\mathbf{w}_1 = -\mathbf{N}_1 \boldsymbol{\xi}_1$  for the well-informed follower UAV. Then, all seven quadrotor models implemented in the MATLAB Simulink. To derive the Nash strategy gain expressions  $\mathbf{K}_i$  and  $\mathbf{N}_i$ , MATLAB's *icare* command is used. The objective functional parameters for (5.17) are selected as

$$\mathbf{Q}_i = \begin{bmatrix} 0.1\mathbf{I}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & 5\mathbf{I}_3 \end{bmatrix}, \quad \mathbf{R}_i = 10\mathbf{I}_4, \quad \gamma = 3,$$

$$c_1 = 0.333, \quad c_2 = 1.5, \quad c_3 = c_4 = 3. \quad (5.57)$$

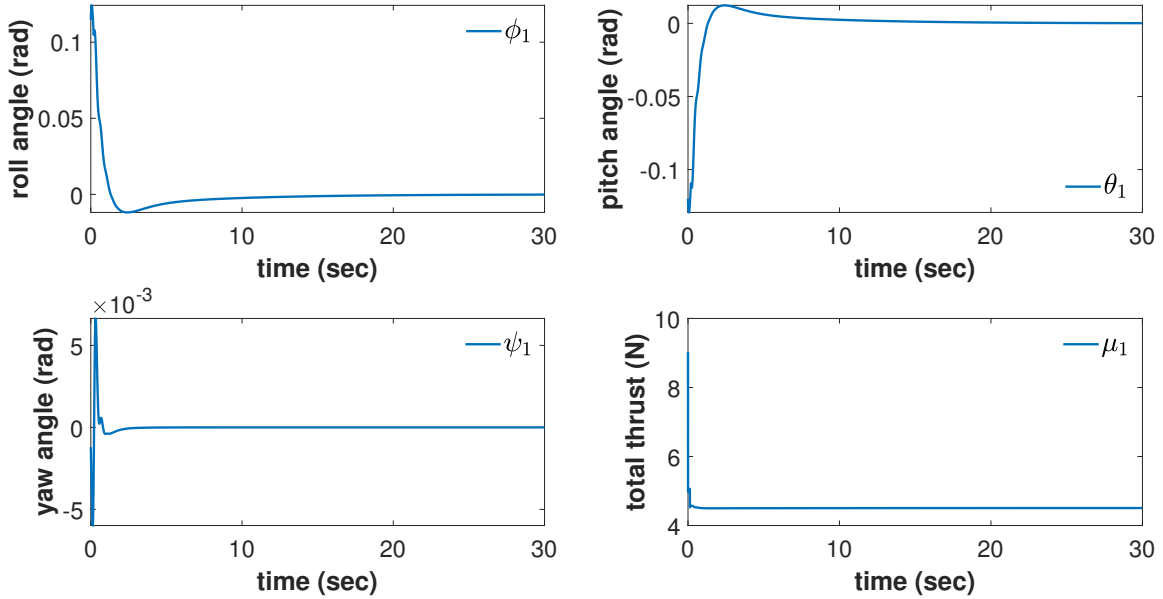


Figure 30: The Euler angles and total thrust of first follower UAV.

The input weighting matrix  $\mathbf{R}_i$  parameters are selected big with respect to the state weighting matrix  $\mathbf{Q}_i$  to make sure that the Euler angles  $\theta$  pitch angle and  $\phi$  roll angle remains in the interval  $[-\pi/6, \pi/6]$ . Note that the condition (5.50) is satisfied for all followers with these parameters.

Since the quadrotor's linear model is obtained around the hover flight condition, the leader's positions are stationary in three-dimensional space. This is shown in the Fig. 28. Additionally, we illustrate attitudes of the quadrotor UAVs in Figures 30-33. The leaders' attitudes at the condition that we linearize the quadrotor's dynamics is hover, and hence their attitude data simply are  $[\theta = 0 \ \phi = 0 \ \psi = c_\psi \ \mu = mg]^T$ .

Realize that the disturbance acting on the well informed first UAV affects not only its stability but also the other uninformed followers. Thence, the uninformed followers' attitudes are nosier with respect to the well-informed UAV and also they become nosier as the nodes of the graph  $G^f$  gets away from the well informed follower node. Nevertheless, the  $\mathcal{H}_\infty$  control method is robust enough to overcome the

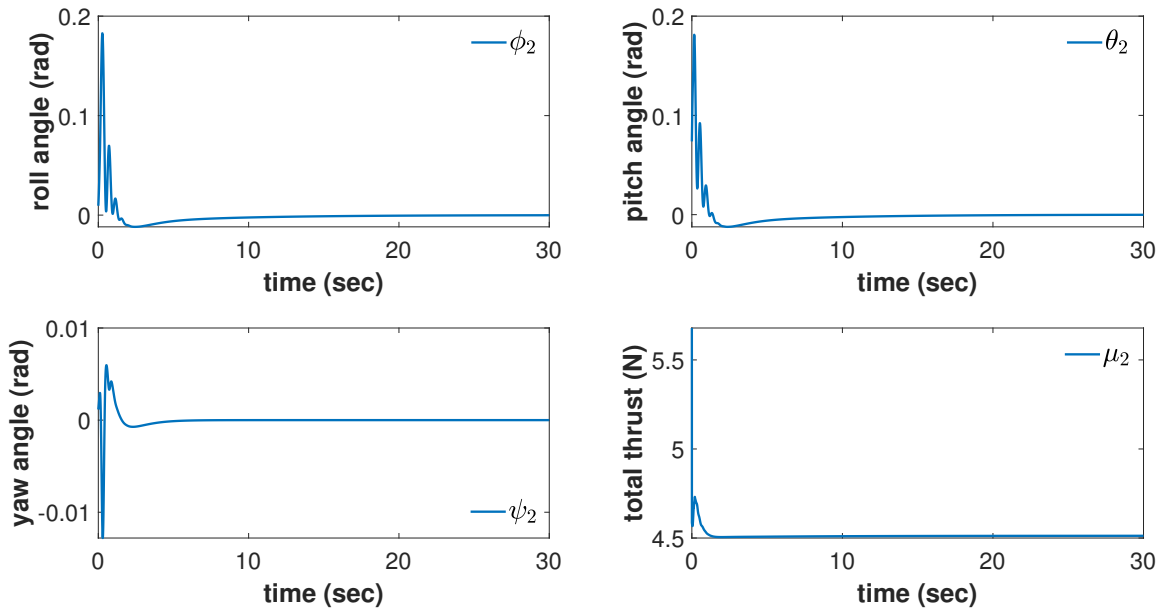


Figure 31: The Euler angles and total thrust of second follower UAV.

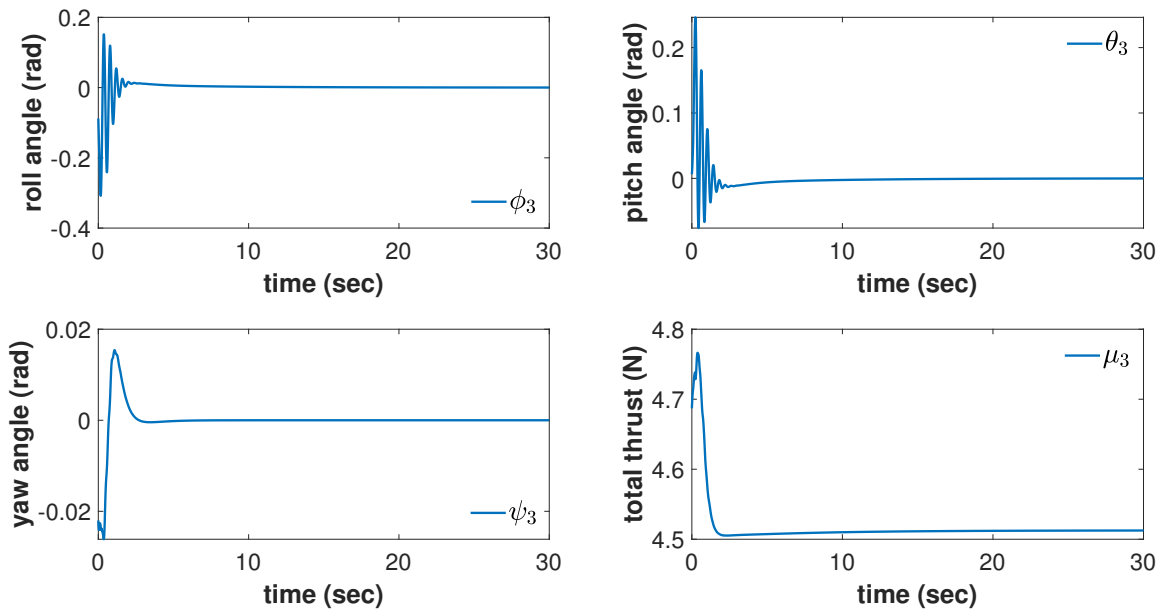


Figure 32: The Euler angles and total thrust of third follower UAV.

disturbance affect as shown in Figures 30-33. Lastly, all of the Euler angle values are converging to zero as expected and the total thrust values are converging to

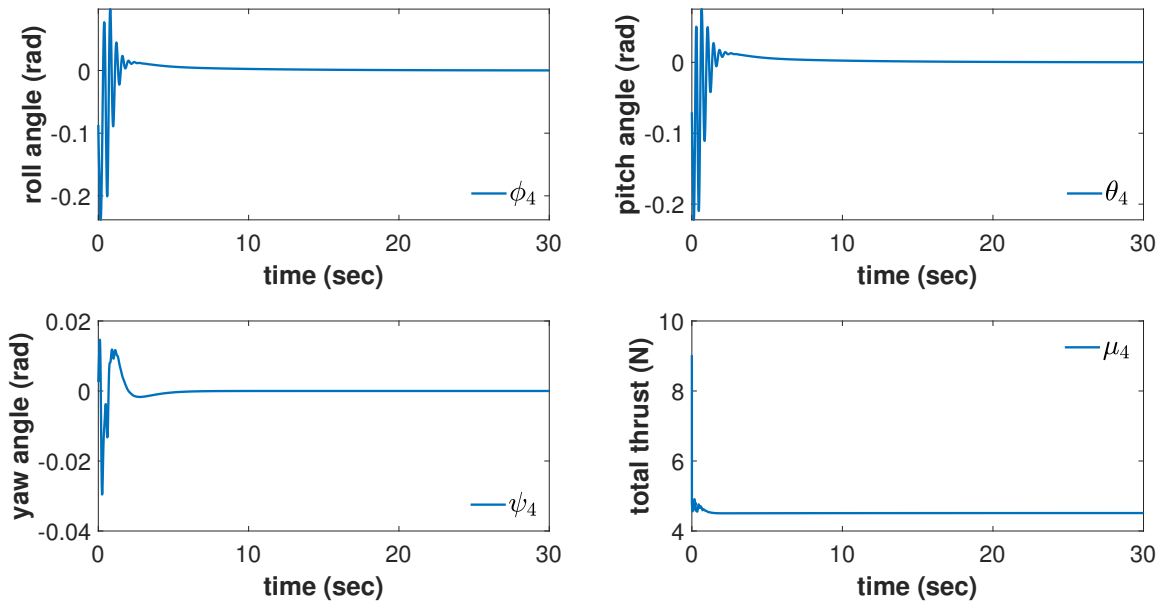


Figure 33: The Euler angles and total thrust of fourth follower UAV.

$mg = 4.5126 \text{ N}$  as can be seen Figures 30-33. This is the value of thrust to make the quadrotors motionless in three-dimensional space.

## CHAPTER 6

### Conclusions and Summary

In the first chapter of this dissertation, a method of distributed backstepping method to have formation flight of multiple quadrotors with distributed time delays is discussed. The proposed algorithms are validated by using Vicon Tracker, AR.Drone 2.0 and a master computer. Through rigorous experimentation and stability analysis, we showed that distributed backstepping control method provides a guaranteed performance for follower agents to track the leader agent with a predetermined position offset. We give the trajectories followed by the single quadrotor under the influence of commensurate delay and by the multiple quadrotors under the influence of distributed delays.

Next, we worked on the game theoretic solution of pursuit-evasion intercept problem when the velocity constraints are imposed on both pursuer and evader. By solving the HJI equation corresponds to the novel non-quadratic functional, we showed that game-optimal velocity trajectories are smooth and satisfies the predetermined boundaries. Using the rigorous Lyapunov analysis, we proved that the PE game ends in a finite-time under certain conditions, which indeed implies that intercept or capture occurs in finite-time. To solve the HJI equation, the IRL method with critic NN structure is used. Consequently, we showed the simulation results of the PE game when the evader adopts both game optimal and sub-optimal velocity policy while the pursuer tracks corresponding game optimal velocity trajectory with the nonlinear backstepping tracker. Simulations showed that when the evader adopts

its game optimal velocity policy, it takes more time to be intercepted by the pursuer compared to the scenario, in which the evader employs a sub-optimal velocity policy.

In the third chapter, we proposed a novel augmented Hamiltonian, and develop a new iterative algorithms based on stationarity conditions of the augmented Hamiltonian to obtain optimal gain solutions for  $H_\infty$  (OPFB) control problem. These gain solutions guarantees both stability and  $L_2$  gain boundedness of an LTI system when the  $H_\infty$  static OPFB control method is employed. Convergence properties of two off-line iterative solution algorithms are given. We showed that solving the Riccati equation iteratively obviates the initial stabilizing gain requirement of Algorithm 1. Then based on Lyapunov iterations, an online off-policy IRL algorithm which is a model-free version of the offline Algorithm 1, is developed to solve the optimal  $H_\infty$  regulator problem by learning the optimal gain solution without requiring system state, control, and disturbance matrices. Lastly, we applied proposed algorithms to the linearized F-16 lateral dynamics at a particular flight condition to verify the correct performance of proposed algorithms.

In the last chapter, we proposed novel local and global objectives for the MLF output containment game for the LTI multi-agent systems. The local  $\mathcal{H}_\infty$  game is solved for the Nash equilibrium strategies where controls are minimizing and disturbances are maximizing players. On the other hand, the graphical game is introduced as a global objective to address mutual interests among the followers. Then, it has been shown that the local and global objectives can be optimized simultaneously under certain conditions. Rigorous  $\mathcal{L}_2$  gain bound and asymptotic stability analyses are provided. The proposed methods are tested by means of the MLF quadrotor UAV game simulations. The results indeed verify the efficacy of the proposed methods.



## REFERENCES

- [1] J. K. Parrish, S. V. Viscido, and D. Grunbaum, “Self-organized fish schools: an examination of emergent properties,” *The biological bulletin*, vol. 202(3), pp. 296–305, 2002.
- [2] D. Mellinger, M. Shomin, N. Michael, and V. Kumar, “Cooperative grasping and transport using multiple quadrotors,” in *Distributed autonomous robotic systems , Heidelberg*. Berlin: Springer, 2013, pp. 545–558.
- [3] T. Kopfstedt, M. Mukai, M. Fujita, and C. Ament, “Control of formations of uavs for surveillance and reconnaissance missions,” *IFAC Proceedings Volumes*, vol. 41(2), pp. 5161–5166, 2008.
- [4] S. Waharte, N. Trigoni, and S. Julier, “Coordinated search with a swarm of uavs,” in *IEEE annual communications society conference on sensor, mesh and ad hoc communications and networks workshops*. IEEE, June 2009, pp. 1–3.
- [5] Y. Kartal, P. Kolaric, V. Lopez, A. Dogan, and F. Lewis, “Backstepping approach for design of pid controller with guaranteed performance for micro-air uav,” *Control Theory and Technology*, vol. 18, no. 1, pp. 19–33, 2020.
- [6] J. M. Daly, Y. Ma, and S. L. Waslander, “Coordinated landing of a quadrotor on a skid-steered ground vehicle in the presence of time delays,” *Autonomous Robots*, vol. 38(2), pp. 179–191, 2015.
- [7] R. Isaacs, *Games of pursuit*. Rand Corporation, 1951.
- [8] H. Modares, M.-B. N. Sistani, and F. L. Lewis, “A policy iteration approach to online optimal control of continuous-time constrained-input systems,” *ISA transactions*, vol. 52, no. 5, pp. 611–621, 2013.

- [9] K. G. Vamvoudakis, D. Vrabie, and F. L. Lewis, “Online adaptive algorithm for optimal control with integral reinforcement learning,” *International Journal of Robust and Nonlinear Control*, vol. 24, no. 17, pp. 2686–2710, 2014.
- [10] J. Gadewadikar, F. L. Lewis, and M. Abu-Khalaf, “Necessary and sufficient conditions for h-infinity static output-feedback control,” *Journal of guidance, control, and dynamics*, vol. 29, no. 4, pp. 915–920, 2006.
- [11] J. Gadewadikar, A. Bhilegaonkar, and F. L. Lewis, “Bounded l2 gain static output feedback: Controller design and implementation on an electromechanical system,” *IEEE Transactions On Industrial Electronics*, vol. 54, no. 5, pp. 2593–2599, 2007.
- [12] A. Al-Tamimi, F. L. Lewis, and M. Abu-Khalaf, “Model-free q-learning designs for linear discrete-time zero-sum games with application to h-infinity control,” *Automatica*, vol. 43, no. 3, pp. 473–481, 2007.
- [13] S. A. A. Rizvi and Z. Lin, “Output feedback q-learning for discrete-time linear zero-sum games with application to the h-infinity control,” *Automatica*, vol. 95, pp. 213–221, 2018.
- [14] Y. Jiang, K. Zhang, J. Wu, C. Zhang, W. Xue, T. Chai, and F. L. Lewis, “H-infinity based minimal energy adaptive control with preset convergence rate,” *IEEE Transactions on Cybernetics*, 2021.
- [15] J. Han, H. Zhang, H. Jiang, and X. Sun, “H-infinity consensus for linear heterogeneous multi-agent systems with state and output feedback control,” *Neurocomputing*, vol. 275, pp. 2635–2644, 2018.
- [16] J. Gadewadikar, F. L. Lewis, K. Subbarao, K. Peng, and B. M. Chen, “H-infinity static output-feedback control for rotorcraft,” *Journal of Intelligent and Robotic Systems*, vol. 54, no. 4, pp. 629–646, 2009.

- [17] Y. Jiang and Z.-P. Jiang, “Computational adaptive optimal control for continuous-time linear systems with completely unknown dynamics,” *Automatica*, vol. 48, no. 10, pp. 2699–2704, 2012.
- [18] H. Modares, F. L. Lewis, and Z.-P. Jiang, “H-infinity 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.
- [19] S. Zuo, Y. Song, F. L. Lewis, and A. Davoudi, “Output containment control of linear heterogeneous multi-agent systems using internal model principle,” *IEEE transactions on cybernetics*, vol. 47, no. 8, pp. 2099–2109, 2017.
- [20] S. Zuo, Y. Song, F. Lewis, and A. Davoudi, “Time-varying output formation containment of general linear homogeneous and heterogeneous multiagent systems,” *IEEE Transactions on Control of Network Systems*, vol. 6, no. 2, pp. 537–548, 2018.
- [21] H. Liu, G. Xie, and L. Wang, “Necessary and sufficient conditions for containment control of networked multi-agent systems,” *Automatica*, vol. 48, no. 7, pp. 1415–1422, 2012.
- [22] K. G. Vamvoudakis, F. L. Lewis, and G. R. Hudas, “Multi-agent differential graphical games: Online adaptive learning solution for synchronization with optimality,” *Automatica*, vol. 48, no. 8, pp. 1598–1611, 2012.
- [23] M. Liu, Y. Wan, V. G. Lopez, F. L. Lewis, G. Hwer, and K. Estabridis, “Differential graphical game with distributed global nash solution,” *IEEE Transactions on Control of Network Systems*, 2021.
- [24] M. I. Abouheaf and F. L. Lewis, “Multi-agent differential graphical games: Nash online adaptive learning solutions,” in *52nd IEEE Conference on Decision and Control*. IEEE, 2013, pp. 5803–5809.

- [25] Z. Qu and M. A. Simaan, “A design of distributed game strategies for networked agents,” *IFAC Proceedings Volumes*, vol. 42, no. 20, pp. 270–275, 2009.
- [26] V. G. Lopez, F. L. Lewis, Y. Wan, M. Liu, G. Hewan, and K. Estabridis, “Stability and robustness analysis of minmax solutions for differential graphical games,” *Automatica*, vol. 121, p. 109177, 2020.
- [27] Y. Kartal, K. Subbarao, N. R. Gans, A. Dogan, and F. Lewis, “Distributed backstepping based control of multiple uav formation flight subject to time delays,” *IET Control Theory & Applications*, vol. 14, no. 12, pp. 1628–1638, 2020.
- [28] I. Maza, A. Ollero, and E. Casado, “‘classification of multi-uav architectures’,” *Handbook of unmanned aerial vehicles*, pp. 953–975, 2015.
- [29] G. Lee and D. Chwa, “Decentralized behavior-based formation control of multiple robots considering obstacle avoidance,” *Intelligent Service Robotics*, vol. 11(1), pp. 127–138, 2018.
- [30] D. Zhou, Z. Wang, and M. Schwager, “Agile coordination and assistive collision avoidance for quadrotor swarms using virtual structures,” *IEEE Transactions on Robotics*, vol. 34(4), pp. 916–923, 2018.
- [31] H. Duan and H. Qiu, “Unmanned aerial vehicle distributed formation rotation control inspired by leader-follower reciprocation of migrant birds,” *IEEE Access*, vol. 6, pp. 23 431–23 443, 2018.
- [32] J. Xi, N. Cai, and Y. Zhong, “Consensus problems for high-order linear time-invariant swarm systems,” *Physica A: Statistical Mechanics and its Applications*, vol. 389(24), pp. 5619–5627, 2010.
- [33] X. Dong, Y. Zhou, and Z. Ren, “‘et al.: ‘time-varying formation control for unmanned aerial vehicles with switching interaction topologies’,” *Control Engineering Practice*, vol. 46, pp. 26–36, 2016.

- [34] J. F. Carvalho, S. Pequito, and A. P. Aguiar, “et al.: ’composability and controllability of structural linear time-invariant systems: Distributed verification’,” *Automatica*, vol. 78, pp. 123–134, 2017.
- [35] C. Q. Ma and J. F. Zhang, “Necessary and sufficient conditions for consensusability of linear multi-agent systems,” *IEEE Transactions on Automatic Control*, vol. 55(5), pp. 1263–1268, 2010.
- [36] W. Rui, D. Xiwang, and L. Qingdong, “et al.: ’adaptive time-varying formation control for high-order lti multi-agent systems’,” *2015 34th Chinese Control Conference (CCC)*, pp. 6998–7003, July 2015.
- [37] F. Baghbani, M. R. Akbarzadeh-T, and M. B. N. Sistani, “Cooperative adaptive fuzzy tracking control for a class of nonlinear multi-agent systems,” in *2017 Joint 17th World Congress of International Fuzzy Systems Association and 9th International Conference on Soft Computing and Intelligent Systems (IFSA-SCIS)*, June 2017, pp. 1–6.
- [38] G. X. Wen, C. P. Chen, and Y. J. Liu, “et al.: ’neural-network-based adaptive leader-following consensus control for second-order non-linear multi-agent systems’,” *IET Control Theory & Applications*, vol. 9(13), pp. 1927–1934, 2015.
- [39] W. He, G. Chen, and Q. L. Han, “et al.: ’network-based leader-following consensus of nonlinear multi-agent systems via distributed impulsive control’,” *Information Sciences*, vol. 380, pp. 145–158, 2017.
- [40] J. Sun and Z. Geng, “Adaptive consensus tracking for linear multi-agent systems with heterogeneous unknown nonlinear dynamics,” *International Journal of Robust and Nonlinear Control*, vol. 26(1), pp. 154–173, 2016.
- [41] M. Defoort, A. Polyakov, and G. Demesure, “et al.: ’leader-follower fixed-time consensus for multi-agent systems with unknown non-linear inherent dynamics’,” *IET Control Theory & Applications*, vol. 9(14), pp. 2165–2170, 2015.

- [42] H. Du, Y. Cheng, and Y. He, “et al.: ‘second-order consensus for nonlinear leader-following multi-agent systems via dynamic output feedback control’,” *International Journal of Robust and Nonlinear Control*, vol. 26(2), pp. 329–344, 2016.
- [43] X. Dong, B. Yu, and Z. Shi, “et al.: ‘time-varying formation control for unmanned aerial vehicles: Theories and applications’,” *IEEE Transactions on Control Systems Technology*, vol. 23(1), pp. 340–348, 2014.
- [44] A. Abdessameud and A. Tayebi, “Formation control of vtol unmanned aerial vehicles with communication delays,” *Automatica*, vol. 47(11), pp. 2383–2394, 2011.
- [45] M. Krstić, I. Kanellakopoulos, and P. V. Kokotović, “Adaptive nonlinear control without overparametrization,” *Systems and Control Letters*, vol. 19(3), pp. 177–185, 1992.
- [46] R. Skjetne, T. I. Fossen, and P. V. Kokotović, “Adaptive maneuvering, with experiments, for a model ship in a marine control laboratory,” *Automatica*, vol. 41(2), pp. 289–298, 2005.
- [47] R. Kristiansen and P. J. Nicklasson, “Satellite attitude control by quaternion-based backstepping,” *Proceedings of the 2005, American Control Conference*, pp. 907–912, June 2005.
- [48] S. S. Pavlichkov, S. N. Dashkovskiy, and C. K. Pang, “Uniform stabilization of nonlinear systems with arbitrary switchings and dynamic uncertainties,” *IEEE Transactions on Automatic Control*, vol. 62(5), pp. 2207–2222, 2016.
- [49] Z. Li and J. Zhao, “Co-design of controllers and a switching policy for nonstrict feedback switched nonlinear systems including first-order feedforward paths,” *IEEE Transactions on Automatic Control*, vol. 64(4), pp. 1753–1760, 2018.

- [50] F. L. Lewis, H. Zhang, K. Hengster-Movric, and A. Das, *Cooperative control of multi-agent systems: optimal and adaptive design approaches*. Springer Science & Business Media, 2013.
- [51] J. Brewer, “Kronecker products and matrix calculus in system theory,” *IEEE Transactions on circuits and systems*, vol. 25(9), pp. 772–781, 1978.
- [52] R. J. Plemmons, “M-matrix characterizations. i—nonsingular m-matrices,” *Linear Algebra and its Applications*, vol. 18(2), pp. 175–188, 1977.
- [53] K. Gu, J. Chen, and V. L. Kharitonov, *Stability of time-delay systems*. Science & Business Media: Springer, 2003.
- [54] Y. Kartal, K. Subbarao, A. Dogan, and F. Lewis, “Optimal game theoretic solution of the pursuit-evasion intercept problem using on-policy reinforcement learning,” *International Journal of Robust and Nonlinear Control*, vol. 31, no. 16, pp. 7886–7903, 2021.
- [55] V. Turetsky and J. Shinar, “Missile guidance laws based on pursuit–evasion game formulations,” *Automatica*, vol. 39, no. 4, pp. 607–618, 2003.
- [56] 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.
- [57] J. R. Marden and J. S. Shamma, “Game theory and control,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, pp. 105–134, 2018.
- [58] R. Isaacs, *Differential games: a mathematical theory with applications to warfare and pursuit, control and optimization*. Courier Corporation, 1999.
- [59] A. E. Bryson, *Applied optimal control: optimization, estimation and control*. CRC Press, 1975.
- [60] F. L. Lewis, D. Vrabie, and V. L. Syrmos, *Optimal control*. John Wiley & Sons, 2012.

- [61] T. Basar and P. Bernhard, *H-infinity optimal control and related minimax design problems: a dynamic game approach*. Springer Science & Business Media, 2008.
- [62] S. Y. Hayoun, M. Weiss, and T. Shima, “A mixed  $l_2/l_\alpha$  differential game approach to pursuit-evasion guidance,” *IEEE Transactions on Aerospace and Electronic Systems*, vol. 52, no. 6, pp. 2775–2788, 2016.
- [63] S. Bhattacharya, T. Başar, and N. Hovakimyan, “A visibility-based pursuit-evasion game with a circular obstacle,” *Journal of Optimization Theory and Applications*, vol. 171, no. 3, pp. 1071–1082, 2016.
- [64] H. Li, D. Liu, and D. Wang, “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.
- [65] Y.-Y. Liu, Z.-S. Wang, and Z. Shi, “Hinf tracking control for linear discrete-time systems via reinforcement learning,” *International Journal of Robust and Nonlinear Control*, vol. 30, no. 1, pp. 282–301, 2020.
- [66] Y. Dong, S. Mingming, and S. Zhaowei, “Satellite proximate interception vector guidance based on differential games,” *Chinese Journal of Aeronautics*, vol. 31, no. 6, pp. 1352–1361, 2018.
- [67] —, “Satellite proximate pursuit-evasion game with different thrust configurations,” *Aerospace Science and Technology*, vol. 99, p. 105715, 2020.
- [68] H. Gong, S. Gong, and J. Li, “Pursuit-evasion game for satellites based on continuous thrust reachable domain,” *IEEE Transactions on Aerospace and Electronic Systems*, 2020.
- [69] A. Jagat and A. J. Sinclair, “Nonlinear control for spacecraft pursuit-evasion game using the state-dependent riccati equation method,” *IEEE Transactions on Aerospace and Electronic Systems*, vol. 53, no. 6, pp. 3032–3042, 2017.



- [70] R. W. Carr, R. G. Cobb, M. Pachter, and S. Pierce, “Solution of a pursuit–evasion game using a near-optimal strategy,” *Journal of Guidance, Control, and Dynamics*, vol. 41, no. 4, pp. 841–850, 2018.
- [71] V. Shaferman and T. Shima, “Cooperative differential games guidance laws for imposing a relative intercept angle,” *Journal of Guidance, Control, and Dynamics*, vol. 40, no. 10, pp. 2465–2480, 2017.
- [72] I. Weintraub, E. Garcia, and M. Pachter, “Optimal guidance strategy for the defense of a non-maneuvrable target in 3-dimensions,” *IET Control Theory & Applications*, vol. 14, no. 11, pp. 1531–1538, 2020.
- [73] M. Abu-Khalaf and F. L. Lewis, “Nearly optimal control laws for nonlinear systems with saturating actuators using a neural network hjb approach,” *Automatica*, vol. 41, no. 5, pp. 779–791, 2005.
- [74] W. M. Haddad and V. Chellaboina, *Nonlinear dynamical systems and control: a Lyapunov-based approach*. Princeton university press, 2011.
- [75] P. Cannarsa and C. Sinestrari, *Semiconcave functions, Hamilton-Jacobi equations, and optimal control*. Springer Science & Business Media, 2004, vol. 58.
- [76] V. G. Lopez, F. L. Lewis, Y. Wan, E. N. Sanchez, and L. Fan, “Solutions for multiagent pursuit-evasion games on communication graphs: Finite-time capture and asymptotic behaviors,” *IEEE Transactions on Automatic Control*, vol. 65, no. 5, pp. 1911–1923, 2019.
- [77] H. Modares and F. L. Lewis, “Optimal tracking control of nonlinear partially-unknown constrained-input systems using integral reinforcement learning,” *Automatica*, vol. 50, no. 7, pp. 1780–1792, 2014.
- [78] Y. Yang, K. G. Vamvoudakis, and H. Modares, “Safe reinforcement learning for dynamical games,” *International Journal of Robust and Nonlinear Control*, vol. 30, no. 9, pp. 3706–3726, 2020.

- [79] H. Jiang, H. Zhang, and X. Xie, “Critic-only adaptive dynamic programming algorithms’ applications to the secure control of cyber–physical systems,” *ISA transactions*, vol. 104, pp. 138–144, 2020.
- [80] A. P. Valadbeigi, A. K. Sedigh, and F. L. Lewis, “H infinity static output-feedback control design for discrete-time systems using reinforcement learning,” *IEEE transactions on neural networks and learning systems*, vol. 31, no. 2, pp. 396–406, 2019.
- [81] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [82] D. P. Bertsekas, *Dynamic programming and optimal control: Vol. 1*. Athena scientific Belmont, 2000.
- [83] 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.
- [84] H. Modares and F. L. Lewis, “Linear quadratic tracking control of partially-unknown continuous-time systems using reinforcement learning,” *IEEE Transactions on Automatic control*, vol. 59, no. 11, pp. 3051–3056, 2014.
- [85] B. Kiumarsi, K. G. Vamvoudakis, H. Modares, and F. L. Lewis, “Optimal and autonomous control using reinforcement learning: A survey,” *IEEE transactions on neural networks and learning systems*, vol. 29, no. 6, pp. 2042–2062, 2017.
- [86] C. Chen, H. Modares, K. Xie, F. L. Lewis, Y. Wan, and S. Xie, “Reinforcement learning-based adaptive optimal exponential tracking control of linear systems with unknown dynamics,” *IEEE Transactions on Automatic Control*, vol. 64, no. 11, pp. 4423–4438, 2019.

- [87] B. Luo, H.-N. Wu, and T. Huang, “Off-policy reinforcement learning for h-infinity control design,” *IEEE transactions on cybernetics*, vol. 45, no. 1, pp. 65–76, 2014.
- [88] M. Abu-Khalaf, F. L. Lewis, and J. Huang, “Neurodynamic programming and zero-sum games for constrained control systems,” *IEEE Transactions on Neural Networks*, vol. 19, no. 7, pp. 1243–1252, 2008.
- [89] H. Zhang, Q. Wei, and D. Liu, “An iterative adaptive dynamic programming method for solving a class of nonlinear zero-sum differential games,” *Automatica*, vol. 47, no. 1, pp. 207–214, 2011.
- [90] Q. Wei, D. Liu, Q. Lin, and R. Song, “Adaptive dynamic programming for discrete-time zero-sum games,” *IEEE transactions on neural networks and learning systems*, vol. 29, no. 4, pp. 957–969, 2017.
- [91] S. Mehraeen, T. Dierks, S. Jagannathan, and M. L. Crow, “Zero-sum two-player game theoretic formulation of affine nonlinear discrete-time systems using neural networks,” *IEEE transactions on cybernetics*, vol. 43, no. 6, pp. 1641–1655, 2012.
- [92] L. M. Zhu, H. Modares, G. O. Peen, F. L. Lewis, and B. Yue, “Adaptive sub-optimal output-feedback control for linear systems using integral reinforcement learning,” *IEEE Transactions on Control Systems Technology*, vol. 23, no. 1, pp. 264–273, 2014.
- [93] F. A. Yaghmaie, S. Gunnarsson, and F. L. Lewis, “Output regulation of unknown linear systems using average cost reinforcement learning,” *Automatica*, vol. 110, p. 108549, 2019.
- [94] 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.

- [95] R. Moghadam and F. L. Lewis, “Output-feedback h-infinity quadratic tracking control of linear systems using reinforcement learning,” *International Journal of Adaptive Control and Signal Processing*, vol. 33, no. 2, pp. 300–314, 2019.
- [96] Q. Jiao, H. Modares, F. L. Lewis, S. Xu, and L. Xie, “Distributed l2-gain output-feedback control of homogeneous and heterogeneous systems,” *Automatica*, vol. 71, pp. 361–368, 2016.
- [97] S. A. Arogeti and F. L. Lewis, “Static output-feedback  $h_\infty$  control design procedures for continuous-time systems with different levels of model knowledge,” *IEEE Transactions on Cybernetics*, 2021.
- [98] J. Gadewadikar, F. L. Lewis, L. Xie, V. Kucera, and M. Abu-Khalaf, “Parameterization of all stabilizing h-infinity static state-feedback gains: application to output-feedback design,” *Automatica*, vol. 43, no. 9, pp. 1597–1604, 2007.
- [99] B. M. Chen, *Robust and H-infinity Control*. Springer Science & Business Media, 2013.
- [100] P. Apkarian, D. Noll, and A. Rondepierre, “Mixed h2 h-infinity control via nonsmooth optimization,” *SIAM Journal on Control and Optimization*, vol. 47, no. 3, pp. 1516–1546, 2008.
- [101] D. Kleinman, “On an iterative technique for riccati equation computations,” *IEEE Transactions on Automatic Control*, vol. 13, no. 1, pp. 114–115, 1968.
- [102] H. W. Knobloch, A. Isidori, and D. Flockerzi, *Topics in control theory*. Birkhäuser, 2012, vol. 22.
- [103] D. Moerder and A. Calise, “Convergence of a numerical algorithm for calculating optimal output feedback gains,” *IEEE Transactions on Automatic Control*, vol. 30, no. 9, pp. 900–903, 1985.
- [104] B. Datta, *Numerical methods for linear control systems*. Academic Press, 2004, vol. 1.

- [105] 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.
- [106] H. Modares, F. L. Lewis, and Z.-P. Jiang, “Optimal output-feedback control of unknown continuous-time linear systems using off-policy reinforcement learning,” *IEEE Transactions on Cybernetics*, vol. 46, no. 11, pp. 2401–2410, 2016.
- [107] C. Li, G. G. Yin, L. Guo, C.-Z. Xu *et al.*, “State observability and observers of linear-time-invariant systems under irregular sampling and sensor limitations,” *IEEE Transactions on Automatic Control*, vol. 56, no. 11, pp. 2639–2654, 2011.
- [108] B. L. Stevens, F. L. Lewis, and E. N. Johnson, *Aircraft control and simulation: dynamics, controls design, and autonomous systems*. John Wiley & Sons, 2015.
- [109] D. Mellinger, M. Shomin, N. Michael, and V. Kumar, *Cooperative grasping and transport using multiple quadrotors*. Springer, 2013.
- [110] T. Kopfstedt, M. Mukai, M. Fujita, and C. Ament, “Control of formations of uavs for surveillance and reconnaissance missions,” *IFAC Proceedings Volumes*, vol. 41, no. 2, pp. 5161–5166, 2008.
- [111] S. Waharte, N. Trigoni, and S. Julier, “Coordinated search with a swarm of uavs,” in *2009 6th IEEE Annual Communications Society Conference on Sensor, Mesh and Ad Hoc Communications and Networks Workshops*. IEEE, 2009, pp. 1–3.
- [112] Y. Kartal, K. Subbarao, A. Dogan, and F. Lewis, “Optimal game theoretic solution of the pursuit-evasion intercept problem using on-policy reinforcement learning,” *International Journal of Robust and Nonlinear Control*, 2021.
- [113] S. B. Sarsilmaz and T. Yucelen, “Distributed control of linear multiagent systems with global and local objectives,” *Systems & Control Letters*, vol. 152, p. 104928, 2021.

- [114] C. Chen, F. L. Lewis, S. Xie, H. Modares, Z. Liu, S. Zuo, and A. Davoudi, “Resilient adaptive and  $h_\infty$  controls of multi-agent systems under sensor and actuator faults,” *Automatica*, vol. 102, pp. 19–26, 2019.
- [115] S. Zuo, F. L. Lewis, and A. Davoudi, “Resilient output containment of heterogeneous cooperative and adversarial multigroup systems,” *IEEE Transactions on Automatic Control*, vol. 65, no. 7, pp. 3104–3111, 2019.
- [116] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [117] D. S. Bernstein, *Matrix mathematics*. Princeton university press, 2009.
- [118] X. Dong and G. Hu, “Time-varying formation tracking for linear multiagent systems with multiple leaders,” *IEEE Transactions on Automatic Control*, vol. 62, no. 7, pp. 3658–3664, 2017.
- [119] F. L. Lewis, H. Zhang, K. Hengster-Movric, and A. Das, *Cooperative control of multi-agent systems: optimal and adaptive design approaches*. Springer Science & Business Media, 2013.
- [120] A. Berman and R. J. Plemmons, *Nonnegative matrices in the mathematical sciences*. SIAM, 1994.
- [121] B. D. McKay, F. E. Oggier, G. F. Royle, N. Sloane, I. M. Wanless, and H. S. Wilf, “Acyclic digraphs and eigenvalues of  $(0, 1)$ -matrices,” *arXiv preprint math/0310423*, 2003.

## BIOGRAPHICAL STATEMENT

Yusuf KARTAL studied Electrical and Electronics Engineering at Bilkent University (Turkey) obtaining the degree of Bachelor's in Science in 2015. Then, he obtained Master in Science degree at Aerospace Engineering Department in University of Texas at Arlington in 2019. He has recently obtained Ph.D. at the same department in University of Texas at Arlington in March 2022. His research interests include multi-robot systems, nonlinear control, game theory, unmanned aerial vehicles, machine learning and distributed control.