

Flocking control against the malicious agent

Chencheng Zhang, Hao Yang, *Senior member, IEEE*, Bin Jiang, *Fellow, IEEE*, and Ming Cao, *Fellow, IEEE*

Abstract—This paper investigates the flocking control of a swarm with a malicious agent that falsifies its controller parameters to cause collision, division, and escape of agents in the swarm. A novel geometric flocking condition is established by designing the configuration of the malicious agent and its neighbors, under which we propose a hierarchical geometric configuration based flocking control method. To help detect the malicious agent, a parameter estimate mechanism is also provided. The proposed method can achieve the flocking control goal and meanwhile contain the malicious agent in the swarm without removing it. Experimental result shows the effectiveness of the theoretical result.

Index Terms—Flocking control; malicious agent; geometric configuration; swarm

I. INTRODUCTION

FLOCKING is a form of collective behavior of plenty of interacting agents with a common group objective under limited environmental information and simple rules. Since Reynolds proposed three heuristic rules: separation, alignment, and cohesion for flocking model in [1], more and more researchers have put effort into the flocking control problem with its applications in multi-agent systems, mobile agents or networks [2], [3]. The main idea of flocking control is to make all agents tend to the same velocity and approach a fixed geometric configuration while preserving the swarm connectivity and avoiding collisions by utilizing artificial intelligence techniques or potential function approaches with local information exchange. In [4], a collection of potential functions are designed for swarms of either single or double integrator agents. Most of these functions are unbounded and are often not appropriate for practice. Therefore, bounded potential functions are investigated by researchers [5], [6]. What's more, many studies appear in the investigation of swarm intelligence for different tasks. For example, Ref. [7] considers the aggregation and formation problem with a discrete-time model. In [8], leader-follower configurations are jointly studied under the model predictive control structure in uncertain environments.

Most of existing results aim at swarm flocking with all agents being healthy and rational. However, agents may suffer from the safety and security issues inevitably in practice. The misbehavior of a swarm appear largely due to three reasons: faults in the physical layer [9], attacks in the cyber layer [10], and abnormal/malicious decisions in the supervisory layer [11], [12].

Under physical faults or cyber attacks, agents may under the appropriate decisions from the supervisory layer. However, malicious decisions refer to the agent's subjectively abnormal and malicious

behavior, which are consequences of either malicious intention or limited cognitive capability of agents. So the malicious agent in the supervisory layer is more difficult to handle. Moreover, since the results on flocking under physical faults are already relatively well developed [13], [14], this paper is devoted to solving flocking problem under abnormal/malicious decisions in the supervisory layer. This is a promising technique with many applications. A typical example is manned-unmanned multiple (air) vehicle swarm where some malicious members may gain control of vehicles to sabotage the mission of the whole swarm [15]. Another example is the well-known Byzantine agents who do not obey the prescribed strategy and update their states arbitrarily to threaten the swarm objective [16]. In real word applications, the control of an Unmanned Aerial Vehicle (UAV) can be taken over by unintended users in a few seconds.

Some effort has been made on control against the malicious agent: For the malicious agent in the cyber layer, the resilient flocking and consensus problems are investigated in [17]- [19]. In these works, although the malicious agent can communicate untruthful information, they still execute the agreed upon decisions. This makes them quite different from the agent with malicious decisions. And these researches consider that the malicious agent can be removed and assume that the network topology remain connected; For the malicious agent in the supervisory layer, Ref. [20] proposes hybrid R -censoring strategies to withstand Byzantine agents and enable cooperative agents to reach consensus. This approach as well as most of other related results merely relies on excluding the malicious agent.

However, to guarantee the completeness of the task in a swarm, the malicious agent is supposed to be safely contained. What's more, the above excluding approaches without considering the motions are not applicable for the networked agents subject to geometric or dynamical constraints such as UAV swarms. To the best of our knowledge, until now almost no result has been reported on *flocking control against malicious decisions of some agent*, let alone the flocking control method that deals with such an agent without excluding it.

Motivated by the above analysis, this paper focuses on the flocking control problem of a swarm in which some agent makes abnormal/malicious decisions in the supervisory layer. Specifically, the malicious agent falsifies its controller parameters, breaks the balance of the attraction or repulsion forces between agents, and thus may lead to collision, division, and escape of agents in the swarm. As a proverb says "one rotten apple could ruin a whole barrel of apples", this paper aims at studying *how the malicious agent affects the whole swarm and how to achieve the flocking control goal without removing the malicious agent from the swarm*. The main contributions of this work are summarized as follows:

A novel geometric flocking condition is established to contain the malicious agent by designing the configuration of the malicious agent and its neighbors, under which the forces acting on the malicious agent from its neighbors reach a balance. We establish a parameter estimate mechanism using filters to help detect the malicious agent with unknown control parameters. Relying on the geometric condition and estimate mechanism, a hierarchical flocking control method is proposed. Such a method consists of the geometrical configuration based control for the neighbors of the malicious agent and the adaptive flocking control for other normal agents. To the best of our knowledge, this is the first attempt to enable a swarm to against the

This work was supported in part by the National Natural Science Foundation of China under Grants 62073165 and 62233009, and in part by the Funding of Postgraduate Research & Practice Innovation Program of Jiangsu Province under Grant KYCX21_0222. (*Corresponding author: Hao Yang*)

C. Zhang is with the College of Automation Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China, Institute of Engineering and Technology (ENTEG), the University of Groningen, the Netherlands (e-mail: zhangchencheng@nuaa.edu.cn).

H. Yang and B. Jiang are with the College of Automation Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing, 211106, China (e-mail: haoyang@nuaa.edu.cn; binjiang@nuaa.edu.cn).

M. Cao is with the Institute of Engineering and Technology (ENTEG), the University of Groningen, the Netherlands (e-mail: m.cao@rug.nl).

agent with malicious decisions and achieve the flocking control goal.

The remainder of the paper is organized as follows: In Section II, preliminaries and model description are given. Sections III provides the malicious agent containment analysis and the flocking control method. The experimental result is presented in Section IV, followed by a conclusion in Section V.

II. PRELIMINARIES

Notations: Let $\mathbf{1}_n$ denote the $n \times 1$ column vector of all ones. Let $|a|_1$ denote the 1-norm and $|a|$ denote the Euclidean-norm of a , respectively. Let $\text{sgn}(a)$ be the signum function of a . Let $\text{diag}(a_1, \dots, a_p)$ be the diagonal matrix with diagonal entries a_1 to a_p . Let $\lambda_{\min}(\cdot)$ denote the minimum eigenvalue of a square real matrix with real eigenvalues. Let \otimes be the Kronecker matrix product.

A. Flocking of a swarm

Consider a swarm of N agents, whose dynamics take the form

$$\begin{cases} \dot{x}_i = v_i, \\ \dot{v}_i = u_i, \end{cases} \quad i \in \mathcal{V}, \quad (1)$$

where $x_i \in \mathbb{R}^m$, $v_i \in \mathbb{R}^m$ and $u_i \in \mathbb{R}^m$ denote the position, the velocity and the control (acceleration) input of agent i for $i \in \mathcal{V}$ with $\mathcal{V} \triangleq \{1, \dots, N\}$. Define $x_{ij} \triangleq x_i - x_j$ as the relative position between agents i and j for $i, j \in \mathcal{V}$. The model (1) can be transformed from a nonlinear flight control system model [21].

The communication topology between agents in swarm (1) is modeled by an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ that consists of a set of vertices \mathcal{V} and a set of edges $\mathcal{E} \triangleq \{(i, j) | i, j \in \mathcal{V}, i \neq j\}$. Vertex $i \in \mathcal{V}$ represents agent i , and edge $(i, j) \in \mathcal{E}$ implies that agents i and j can interact with each other and are unordered. An *undirected path* between vertices i and j is a sequence of unordered edges, $(i, k_1), (k_1, k_2), \dots, (k_l, j)$ with distinct vertices k_p , $p = 1, 2, \dots, l$. If there exists an undirected path between vertices i and j , the two vertices are said to be *connected*; otherwise, they are *unconnected*. An undirected graph is called *connected* if any two distinct vertices in the graph are connected. The Laplacian matrix of graph \mathcal{G} is denoted by L . Define R as the *sensing radius* of each agent, which indicates that two agents can interact only if distance between them is smaller than R , i.e., if $0 < |x_{ij}| < R$, then $(i, j) \in \mathcal{E}$; otherwise, $(i, j) \notin \mathcal{E}$. Agent j is called a *neighbor* of agent i if $(i, j) \in \mathcal{E}$. Define $\mathcal{N}(i) \triangleq \{j \in \mathcal{V} : (i, j) \in \mathcal{E}, i \neq j\}$ as the set of neighbors of agent i in \mathcal{G} . Note that the following study can be applied to the case that the communication topology is considered static as well.

The flocking control objective is to *make the whole swarm tend to a common speed and approach a fixed configuration without collision*, i.e., $\lim_{t \rightarrow \infty} \dot{x}_{ij} = \lim_{t \rightarrow \infty} v_i - v_j = 0$, $\forall i, j \in \mathcal{V}$; $0 < |x_{ik}(t)| < R$, $t \geq 0$, $\forall i \in \mathcal{V}$, $k \in \mathcal{N}(i)$. A conventional flocking control law is designed as [2]

$$u_i = - \sum_{j \in \mathcal{N}(i)} (v_i - v_j) - \sum_{j \in \mathcal{N}(i)} \nabla_{x_i} V_{ij}(|x_{ij}|), \quad i \in \mathcal{V} \quad (2)$$

where the first term corresponds to the desired velocity alignment, and the second term is the gradient of a potential function V_{ij} . Note that many existing potential functions with different forms can be applied here in the normal case, for example, the bounded potential function proposed in [6]

$$V_{ij}(|x_{ij}|) \triangleq \underbrace{\frac{R^2 - |x_{ij}|^2}{|x_{ij}|^2 + \frac{R^2}{E}}}_{V_{rij}} + \underbrace{\frac{|x_{ij}|^2}{R^2 - |x_{ij}|^2 + \frac{R^2}{E}}}_{V_{aij}}, \quad 0 \leq |x_{ij}| \leq R \quad (3)$$

where E is a positive constant. V_{ij} satisfies the following properties

- $V_{ij}(|x_{ij}|) = E$ when $|x_{ij}| = 0$ or $|x_{ij}| = R$;
- $\frac{\partial V_{ij}(|x_{ij}|)}{\partial(|x_{ij}|)} < 0$ when $|x_{ij}| \in (0, \delta)$ and $\frac{\partial V_{ij}(|x_{ij}|)}{\partial(|x_{ij}|)} > 0$ when $|x_{ij}| \in (\delta, R)$, where $\delta \triangleq \frac{\sqrt{2}R}{2}$.

Physically, the potential can be divided into $V_{ij} \triangleq V_{aij} + V_{rij}$ where V_{aij} and V_{rij} can be viewed as potentials of attraction and repulsion of agent i with respect to agent j , respectively. Obviously, V_{ij} reaches its minimum when $|x_{ij}| = \delta$. In the unique distance δ , it holds that $\nabla_{x_i} V_{aij}(\delta) + \nabla_{x_i} V_{rij}(\delta) = 0$. In normal case, one can choose $E > \bar{Q} \triangleq \frac{1}{2} \sum_{i \in \mathcal{V}} v_i^T(0) v_i(0) + \frac{N(N-1)}{2} \max_{i, j \in \mathcal{V}} \{\bar{V}_{ij}(|x_{ij}(0)|)\}$ where $\bar{V}_{ij}(|x_{ij}|) \triangleq \frac{R^2 - |x_{ij}|^2}{|x_{ij}|^2} + \frac{|x_{ij}|^2}{R^2 - |x_{ij}|^2}$. This makes the potential between any two agents sufficiently large when the distance between them is equal to 0 or R , and thus avoids the collision while preserving the connectivity [6]. In the sequel, E will be chosen sufficiently large (i.e., larger than the initial energy functions built in the following sections) to avoid the collision and preserve the connectivity when applying the potential function V_{ij} in the control design.

B. Malicious agent

Consider a malicious agent, denoted as $i_f \in \mathcal{V}$, who intentionally falsifies controller parameters such that

$$u_{i_f} = -k_v \sum_{j \in \mathcal{N}(i_f)} (v_{i_f} - v_j) - \sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} \tilde{V}_{i_f j}(|x_{i_f j}|) \quad (4)$$

where

$$\tilde{V}_{i_f j}(|x_{i_f j}|) \triangleq k_a V_{ai_f j}(|x_{i_f j}|) + k_r V_{ri_f j}(|x_{i_f j}|) \quad (5)$$

We provide some insights on these parameters:

- k_a : This parameter represents the strength of the attractive force on agent i_f which is inverted for $k_a < 0$, completely lost for $k_a = 0$, partially lost for $0 < k_a < 1$, and strengthened for $k_a > 1$.
- k_r : This parameter represents the strength of the repulsive force on agent i_f which is inverted for $k_r < 0$, completely lost for $k_r = 0$, partially lost for $0 < k_r < 1$, and strengthened for $k_r > 1$.
- $k_v < 1$: This parameter represents the efficacy for the velocity consensus of agent i_f which is inverted for $k_v < 0$, completely lost for $k_v = 0$, and partially lost for $0 < k_v < 1$.

Compared with the normal controller in (2), the attraction/repulsion effort acting on agent i_f from each of its neighbors is out of balance under the distance δ . The resultant force of agent i_f is decided by the combination of these three parameters. When $k_v = k_a = k_r = 1$, the malicious agent is degenerated into a normal one.

Specifically, there are two circumstances that can cause serious influence to the swarm: (1) When $k_v \leq 0$, $k_a \leq 0$ and $k_r \gg 1$, the malicious agent i_f tries its best to *run away* from the agents around it and may finally escape from the swarm; (2) When $k_v \leq 0$, $k_r \leq 0$ and $k_a \gg 1$, agent i_f tries its best to *collide* with the agents around it.

Assumption 1 : $|k_v| \leq \bar{k}_v$, $|k_a| \leq \bar{k}_a$, $|k_r| \leq \bar{k}_r$ for $\bar{k}_v, \bar{k}_a, \bar{k}_r > 0$. \square

This assumption means that these parameters are bounded and this is helpful to design the bounds of potential functions. Such an assumption is not required if the unbounded potential functions instead of the bounded ones are applied in this research.

In the following, the definition of containing a malicious agent is presented.

Definition 1 : The malicious agent i_f is said to be *contained* if $\dot{v}_{i_f} = u_{i_f} = 0$ and $|x_{i_f j}| = \bar{\delta}_{i_f j}$ where $0 < \bar{\delta}_{i_f j} < R$ is a designable expected distance between agent i_f and its neighbor $j \in \mathcal{N}(i_f)$. \square

C. Problem formulation

Define a undirected graph $\mathcal{G}' \triangleq (\mathcal{V}', \mathcal{E}')$ consisting of the set of vertices $\mathcal{V}' \triangleq \mathcal{V} - \{i_f\}$ and the set of edges $\mathcal{E}' \triangleq \{(i, j) | i, j \in \mathcal{V}', |x_{ij}| < R, i \neq j\}$. Define $V_i \triangleq \sum_{j \in \mathcal{N}(i)} V_{ij}$ for $i \in \mathcal{V}, j \in \mathcal{N}(i)$ and $\tilde{V}_{i_f} \triangleq \sum_{j \in \mathcal{N}(i_f)} \tilde{V}_{i_f j}$.

Assumption 2 : The initial graph \mathcal{G}' is connected. \square

Assumption 2 guarantees that the information among all the normal agents in the swarm can be transmitted at the initial time. Similar classical assumption on initial graph can be found in many flocking control researches such as [3], [5]. Based on this assumption, the problem to be solved is formulated as follows.

Problem \mathcal{F} : Consider the swarm (1) satisfying Assumptions 1-2 with a malicious agent $i_f \in \mathcal{V}$ under controller (4)-(5). Design $u_i, i \in \mathcal{V}'$ such that

① $\lim_{t \rightarrow \infty} (v_i - v_j) = 0, \forall i, j \in \mathcal{V}$, i.e., all the agents tend to a same velocity;

② The swarm (1) asymptotically converges to a fixed geometric configuration, under which

- $u_{i_f} = 0$ and $|x_{i_f j}| = \bar{\delta}_{i_f j}, \forall j \in \mathcal{N}(i_f)$ where $0 < \bar{\delta}_{i_f j} < R$, i.e., the malicious agent i_f is contained.
- $|x_{ij}| = \bar{\delta}_{ij}, \forall i \in \mathcal{V}', j \in \mathcal{N}(i)$ where $0 < \bar{\delta}_{ij} < R$, i.e., the normal agents are connected with their neighbors;

③ $|x_{ij}(t)| \neq 0$ for $t \geq 0, \forall i, j \in \mathcal{V}$ and $i \neq j$, i.e., no collision occurs. \square

III. MAIN RESULT

A. Malicious agent containment analysis

We first establish a distributed geometric condition under which the malicious agent can still be contained in the swarm. Such a condition will be the basis for the flocking control design.

Lemma 1 : Consider the swarm (1) with malicious agent $i_f \in \mathcal{V}$ under controller (4)-(5). Suppose that $v_{i_f} - v_j = 0, \forall j \in \mathcal{N}(i_f)$. If

$$|x_{i_f j}| = \bar{\delta}, \forall j \in \mathcal{N}(i_f), \quad (6)$$

$$\sum_{j \in \mathcal{N}(i_f)} x_{i_f j} = 0 \quad (7)$$

Proof : As $\tilde{V}_{i_f j}$ defined in (5) is symmetric with respect to $x_{i_f j}$, it holds that

$$\begin{aligned} \sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} \tilde{V}_{i_f j}(|x_{i_f j}|) &= \sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f j}} \tilde{V}_{i_f j}(|x_{i_f j}|) \\ &= \sum_{j \in \mathcal{N}(i_f)} \frac{\partial \tilde{V}_{i_f j}(|x_{i_f j}|)}{\partial |x_{i_f j}|} \cdot \frac{\partial |x_{i_f j}|}{\partial x_{i_f j}} \end{aligned}$$

It yields from the definition of Euclidean norm that for $j \in \mathcal{N}(i_f)$

$$\begin{aligned} \frac{\partial |x_{i_f j}|}{\partial x_{i_f j}} &= \frac{\partial (x_{i_f j}^T x_{i_f j})^{\frac{1}{2}}}{\partial x_{i_f j}} = \frac{1}{2} (x_{i_f j}^T x_{i_f j})^{-\frac{1}{2}} \cdot \frac{\partial (x_{i_f j}^T x_{i_f j})}{\partial x_{i_f j}} \\ &= \frac{2x_{i_f j}}{2(x_{i_f j}^T x_{i_f j})^{\frac{1}{2}}} = \frac{x_{i_f j}}{|x_{i_f j}|} \end{aligned}$$

Thus, one can obtain that

$$\sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} \tilde{V}_{i_f j}(|x_{i_f j}|) = \sum_{j \in \mathcal{N}(i_f)} \frac{\partial \tilde{V}_{i_f j}(|x_{i_f j}|)}{\partial |x_{i_f j}|} \cdot \frac{x_{i_f j}}{|x_{i_f j}|} \quad (8)$$

Define s as the number of agents in $\mathcal{N}(i_f)$. Condition (6) indicates that $|x_{i_f j_1}| = |x_{i_f j_2}| = \dots = |x_{i_f j_s}|$ for $j_1, j_2, \dots, j_s \in \mathcal{N}(i_f)$. Therefore, it holds $\tilde{V}_{i_f j_1}(|x_{i_f j_1}|) = \tilde{V}_{i_f j_2}(|x_{i_f j_2}|) = \dots = \tilde{V}_{i_f j_s}(|x_{i_f j_s}|)$. It yields from (8) that

$$\begin{aligned} &\sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} \tilde{V}_{i_f j}(|x_{i_f j}|) \\ &= s \frac{\partial \tilde{V}_{i_f j}(|x_{i_f j}|)}{\partial |x_{i_f j}|} \Big|_{|x_{i_f j}|=\bar{\delta}} \cdot \frac{(x_{i_f j_1} + x_{i_f j_2} + \dots + x_{i_f j_s})}{\bar{\delta}} \end{aligned}$$

According to condition (7), one further has

$$\begin{aligned} &\sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} \tilde{V}_{i_f j}(|x_{i_f j}|) \\ &= s \frac{\partial \tilde{V}_{i_f j}(|x_{i_f j}|)}{\partial |x_{i_f j}|} \Big|_{|x_{i_f j}|=\bar{\delta}} \cdot \frac{\sum_{j \in \mathcal{N}(i_f)} x_{i_f j}}{\bar{\delta}} = 0 \end{aligned}$$

Suppose that $v_{i_f} - v_j = 0$ for $j \in \mathcal{N}(i_f)$. According to the malicious controller (4)-(5) of i_f , under conditions (6)-(7), one can deduce that

$$\begin{aligned} \dot{v}_{i_f} = u_{i_f} &= - \sum_{j \in \mathcal{N}(i_f)} k_v (v_{i_f} - v_j) - \sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} \tilde{V}_{i_f j}(|x_{i_f j}|) \\ &= 0 \end{aligned} \quad (9)$$

This completes the proof. \square

Remark 1 : Conditions (6)-(7) provide a desired geometrical configuration that is a regular polygon with the malicious agent being the center and its neighbors being vertexes. In this case, the total potential gradient of the malicious agent with respect to its neighbors is restricted to be 0 and their distances are also fixed. Physically, this means that the forces acted on the malicious agent from all its neighbors reach a balance such that the malicious agent can still be contained in the swarm. An example of the desired configuration satisfying (6)-(7) is presented in Fig. 1, where the malicious agent is surrounded by three neighbors. \square

Remark 2 : Conditions (6)-(7) require the number of agents in $\mathcal{N}(i_f)$, $s \geq 2$. This is because when agent i_f has at least two neighbors, there exist expected extreme points of the total potential \tilde{V}_{i_f} such that $\nabla_{x_{i_f}} \tilde{V}_{i_f}$ can be 0. Then agent i_f 's malicious behavior can be contained. Provided that $s = 1$, \tilde{V}_{i_f} is only related to $|x_{i_f j}|$ for $j \in \mathcal{N}(i_f)$. According to the malicious controller (4)-(5), \tilde{V}_{i_f} tries to reach its minimum. However, as \tilde{V}_{i_f} reaches its minimum, $|x_{i_f j}|$ reaches an unexpected or even dangerous distance, for example, $|x_{i_f j}| = 0$ when $k_r = 0$ and $k_a \neq 0$ in (7). No expected extreme point can be found since \tilde{V}_{i_f} monotonically increases with respect to $|x_{i_f j}|$. Once $|x_{i_f j}| \neq 0$, it holds $\nabla_{x_{i_f}} \tilde{V}_{i_f} > 0$ which leads to the acceleration of agent i_f . Thus, the malicious agent can never be contained. \square

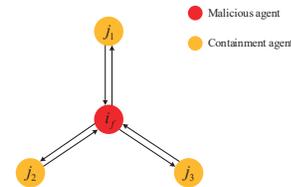


Fig. 1. An illustration of the desired configuration to contain the malicious agent.

B. Hierarchical geometric configuration based flocking control

In this subsection, a hierarchical geometric configuration based flocking control method is proposed to solve problem \mathcal{F} . The control

architecture is shown in Fig. 2, where the malicious agent is in Layer 1, all its neighbors are in Layer 2, and other agents in the swarm are in Layer 3. An important feature of such an architecture is that the agents in Layer 2 do not utilize the information of agents in Layer 3. This feature makes agents in Layer 2 form the desired configuration as in Lemma 1 more conveniently to contain the malicious one. In this case, the agents in Layers 2 and 3 can be viewed as *leaders* and *followers*, respectively. Define \mathcal{V}_l as the set of agents in Layer 2, \mathcal{V}_f as the set of agents in Layer 3, and $\mathcal{V}_g \triangleq \{i_f\} + \mathcal{N}(i_f)$ as the set of agents in Layers 1 and 2 as shown in Fig. 2. Next we shall design controllers for agents in Layer 2 and Layer 3, respectively.

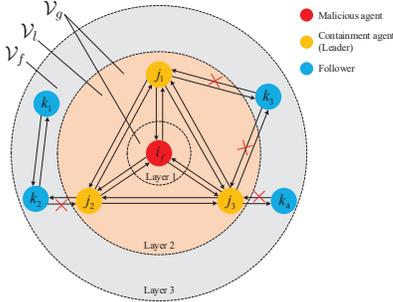


Fig. 2. An illustration of the hierarchical control architecture.

Before giving the main result, the following assumptions are made.

Assumption 3 : At the initial time, there are at least two neighbors of the malicious agent. \square

Assumption 4 : At the initial time, any two agents in $\mathcal{N}(i_f)$ are neighbors. \square

Assumption 3 is a condition on the number of the malicious agent's neighbors under the geometric configuration method, which has been explained in Remark 2. Assumption 4 means that all neighbors of the malicious agent can interact with each other at the initial time. Such an assumption is needed to resist the influence of the malicious agent by its neighbors jointly under a distributed control structure, and will be explained in details in Remarks 3 and 4.

For convenience, rewrite the dynamics of the malicious agent i_f

$$\dot{v}_{i_f} = -C_{i_f} k \quad (10)$$

where $k \triangleq (k_v, k_a, k_r)^T$ is the vector of the unknown parameters and $C_{i_f} \triangleq \left(\sum_{j \in \mathcal{N}(i_f)} (v_{i_f} - v_j), \sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} V_{aifj}(|x_{ifj}|), \sum_{j \in \mathcal{N}(i_f)} \nabla_{x_{i_f}} V_{rifj}(|x_{ifj}|) \right)$.

In order to track the unknown parameters, filter v_{i_f} and C_{i_f} in (10) by low-pass first-order filters, one has

$$\dot{v}_{i_f}^F = -a v_{i_f}^F + a v_{i_f}, \quad v_{i_f}^F(0) = v_{i_f}(0) \quad (11)$$

$$\dot{C}_{i_f}^F = -a C_{i_f}^F + C_{i_f}, \quad C_{i_f}^F(0) = 0 \quad (12)$$

where $a > 0$ is the scalar filter gain. $v_{i_f}^F$ and $C_{i_f}^F$ are the filtered v_{i_f} and C_{i_f} , respectively. They can be obtained by the above stable linear filter equations (11)-(12). And it holds that $\dot{v}_{i_f}^F = -a C_{i_f}^F k$. This together with (11) yields

$$-v_{i_f}^F + v_{i_f} = -C_{i_f}^F k \quad (13)$$

Define $\hat{k} \triangleq (\hat{k}_v, \hat{k}_a, \hat{k}_r)^T$ as the estimate of k . Design the adaptive update law of the estimate as follows

$$\begin{aligned} \dot{\hat{k}} = & -\Gamma_k C_{i_f}^T \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f}) \\ & - \Gamma_k (C_{i_f}^F)^T (C_{i_f}^F \hat{k} + v_{i_f} - v_{i_f}^F) \end{aligned} \quad (14)$$

where Γ_k is the positive-definite gain matrix.

Based on Conditions (6)-(7) of Lemma 1, let $\sum_{j \in \mathcal{N}(i_f)} x_{ifj}^* = 0$ and $|x_{ifj}^*| = \bar{\delta} < R/2$ where x_{ifj}^* denotes the desired displacement between agents i_f and $j \in \mathcal{N}(i_f)$. Note that $x_{jk} = x_{ji_f} - x_{ki_f}$ and $x_{jk}^* = x_{ji_f}^* - x_{ki_f}^*, \forall j, k \in \mathcal{N}(i_f)$. Design the controller of agent $j \in \mathcal{N}(i_f)$ as follows

$$\begin{aligned} u_j = & -\kappa_v \sum_{p \in \mathcal{N}(j) \cap \mathcal{V}_l} (v_j - v_p) - \kappa_x \sum_{p \in \mathcal{N}(j) \cap \mathcal{V}_l} \nabla_{x_j} \hat{V}_{jp}(x_{jp}) \\ & - C_{i_f} \hat{k} \end{aligned} \quad (15)$$

where constants $\kappa_v, \kappa_x \geq 1$. The non-negative potential function $\hat{V}_{ij}(x_{ij})$ satisfies the following properties that

- 1) \hat{V}_{ij} attains its unique minimum and $\frac{\partial \hat{V}_{ij}}{\partial |x_{ij} - x_{ij}^*|} = 0$ when $x_{ij} = x_{ij}^*$;
- 2) $\hat{V}_{ij} > \bar{H}$ when $|x_{ij}| = 0$ and $|x_{ij}| = R$ where \bar{H} is a designable positive constant.

To solve Problem \mathcal{F} , \bar{H} is chosen as follows

$$\begin{aligned} \bar{H} = & \sum_{j \in \mathcal{N}(i_f)} \left(\kappa_x \hat{V}'_{ji_f}(0) + \frac{\kappa_x}{2} \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} \hat{V}'_{ji}(0) + \right. \\ & \left. \frac{1}{2} (v_j(0) - v_{i_f}(0))^T (v_j(0) - v_{i_f}(0)) \right) + \frac{1}{2} \lambda_{\max}(\Gamma_k^{-1}) \times \\ & \left((\bar{k}_v + \hat{k}_v(0))^2 + (\bar{k}_a + \hat{k}_a(0))^2 + (\bar{k}_r + \hat{k}_r(0))^2 \right) \end{aligned} \quad (16)$$

where $\hat{V}'_{ij}(0) \triangleq \frac{|x_{ij}(0) - x_{ij}^*|^2}{R - |x_{ij}(0)|} + \frac{|x_{ij}(0) - x_{ij}^*|^2}{|x_{ij}(0)|}$.

Condition 1) shows that the potential between two agents is minimized when their displacement is equal to the desired one, which makes the two agents approach to the desired configuration. Condition 2) means that the potential would become sufficiently large when the two agents tend to collide or escape, and thus guarantees that no collision happens and no edge is lost. One example of such a potential function is as follows

$$\hat{V}_{ij}(x_{ij}) \triangleq \frac{|x_{ij} - x_{ij}^*|^2}{R - |x_{ij}| + \frac{(R - \delta_{ij})^2}{H + \iota}} + \frac{|x_{ij} - x_{ij}^*|^2}{|x_{ij}| + \frac{\delta_{ij}^2}{H + \iota}}$$

for $0 \leq |x_{ij}| \leq R$, where ι is a positive constant and $\delta_{ij} \triangleq |x_{ij}^*|, 0 < |x_{ij}^*| < R$.

Remark 3 : Note that in controller (15), agent $j \in \mathcal{N}(i_f)$ only utilizes the information in \mathcal{V}_g (Layers 1-2) rather than information in \mathcal{V}_f (Layer 3). In the last term of controller (15), the estimate \hat{k} of the unknown parameter k with adaptive updating law (14) is utilized. And controller (15) requires the state information among all the neighbors of the malicious agent. This local information exchange is required since all neighbors need to jointly resist the influence of the malicious agent. As will be shown in Theorem 1, under Assumption 4, this local information exchange is always available. We shall explain this setting later in Remark 4. \square

Now design a distributed adaptive controller for agent $k \in \mathcal{V}_f$ as

$$u_k = - \sum_{p \in \mathcal{N}(k)} \alpha_{kp} \text{sgn}(v_k - v_p) - \sum_{p \in \mathcal{N}(k)} \nabla_{x_k} V_{kp}(|x_{kp}|) \quad (17)$$

$$\dot{\alpha}_{kp} = \gamma_{kp} |v_k - v_p|, \quad p \in \mathcal{N}(k)$$

where α_{kp} is a varying gain with initial values $\alpha_{kp}(0) \geq 0$ and V_{kp} is defined in (2). γ_{kp} is a positive constant and $\gamma_{kp} = \gamma_{pk}$.

Theorem 1 : Consider the swarm (1) satisfying Assumptions 1-4 with malicious agent $i_f \in \mathcal{V}$ under controller (4)-(5). Problem \mathcal{F} is solved by applying controller (15) along with parameter estimate

update law (14) to agents in \mathcal{V}_l and controller (17) to agents in \mathcal{V}_f . \square

Before moving on, the following concepts of directed graph theory and a lemma are given that will be used to prove Theorem 1.

A *directed graph* $\hat{\mathcal{G}}$ consists of a pair $(\hat{\mathcal{V}}, \hat{\mathcal{E}})$, where $\hat{\mathcal{V}} \triangleq \{1, \dots, p\}$ is a set of vertices and $\hat{\mathcal{E}} \triangleq \{(j, k) | j, k \in \hat{\mathcal{V}}, j \neq k\}$ is a set of ordered pairs of vertices. An edge (j, k) denotes that vertex k can obtain and utilize information from vertex j , but not necessarily vice versa. A *directed path* from vertex j to k is a sequence of edges denoted by $(j, i_1), (i_1, i_2), \dots, (i_r, k)$ with $i_l \in \hat{\mathcal{V}}, l \in \{1, \dots, r\}$.

Lemma 2 [23] : For an undirected connected graph with the Laplacian matrix $L \in \mathfrak{R}^{n \times n}$, given $B \triangleq \text{diag}(b_1, \dots, b_n)$ where $b_i \geq 0, i = 1, \dots, n$, if there exists $b_i > 0$, then the matrix $E = L + B$ is symmetric positive definite. \square

Proof of Theorem 1 : The proof of Theorem 1 is divided into two parts: In *Part A*, we prove that $\lim_{t \rightarrow \infty} (v_i - v_j) = 0, \lim_{t \rightarrow \infty} x_{ji} = x_{ji}^*$ and $0 < |x_{ij}(t)| < R$ for $t \geq 0, \forall i, j \in \mathcal{V}_g$. In *Part B*, we prove that $\lim_{t \rightarrow \infty} (v_a - v_b) = 0$ and $0 < |x_{ab}(t)| < R$ for $t \geq 0, \forall a, b \in \mathcal{V}$. $\lim_{t \rightarrow \infty} \nabla_{x_p} V_p = 0, \forall p \in \mathcal{V}_f$.

Part A. The behavior of agents in Layers 1-2 is considered and the proof is given by analyzing the error velocity and position information between the malicious agent and its neighbors. Denote $x_f \triangleq (x_{i_f}^T, x_{j_1}^T, \dots, x_{j_s}^T)^T, v_f \triangleq (v_{i_f}^T, v_{j_1}^T, \dots, v_{j_s}^T)^T$ for $j_k \in \mathcal{N}(i_f), k \in \{1, \dots, s\}$. Define an energy function $H(x_f, v_f)$ as

$$\begin{aligned} H(x_f, v_f) &\triangleq \kappa_x \sum_{j \in \mathcal{N}(i_f)} \hat{V}_{ji_f} + \frac{\kappa_x}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} \hat{V}_{ji} \\ &+ \frac{1}{2} \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f})^T (v_j - v_{i_f}) \\ &+ \frac{1}{2} \tilde{k}^T \Gamma_k^{-1} \tilde{k} \end{aligned} \quad (18)$$

where $\tilde{k} \triangleq (\tilde{k}_v, \tilde{k}_a, \tilde{k}_r)^T$ with $\tilde{k}_v \triangleq k_v - \hat{k}_v, \tilde{k}_a \triangleq k_a - \hat{k}_a, \tilde{k}_r \triangleq k_r - \hat{k}_r$. Note that $\hat{V}_{ji} = \hat{x}_{ji} \nabla_{x_{ji}} \hat{V}_{ji}$. The time derivative of $H(x_f, v_f)$ is

$$\begin{aligned} \dot{H}(x_f, v_f) &= \kappa_x \sum_{j \in \mathcal{N}(i_f)} (v_j^T - v_{i_f}^T) \nabla_{x_j} \hat{V}_{ji_f} \\ &+ \frac{\kappa_x}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} (v_j^T - v_i^T) \nabla_{x_j} \hat{V}_{ji} \\ &+ \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f})^T (\dot{v}_j - \dot{v}_{i_f}) - \tilde{k}^T \Gamma_k^{-1} \dot{\tilde{k}} \end{aligned}$$

Applying the filters (11)-(12), the estimator (14) and controller (15), one has

$$\begin{aligned} \dot{H}(x_f, v_f) &= \kappa_x \sum_{j \in \mathcal{N}(i_f)} \left((v_j^T - v_{i_f}^T) \nabla_{x_j} \hat{V}_{ji_f} + \frac{1}{2} \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} (v_j^T - v_i^T) \right. \\ &\quad \times \nabla_{x_j} \hat{V}_{ji} \left. \right) + \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f})^T \left(-\kappa_v (v_j - v_{i_f}) \right. \\ &\quad - \kappa_v \sum_{p \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} (v_j - v_p) - \kappa_x \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} \nabla_{x_j} \hat{V}_{ji} \\ &\quad \left. - \kappa_x \nabla_{x_j} \hat{V}_{ji_f} \right) - \tilde{k}^T (C_{i_f}^F)^T C_{i_f}^F \tilde{k} \end{aligned} \quad (19)$$

It follows from the fact that $x_{ji} = -x_{ij}$ and $x_{ji} - x_{ji}^* = -(x_{ij} - x_{ij}^*)$

that $\frac{\partial \hat{V}_{ji}}{\partial x_{ji}} = \frac{\partial \hat{V}_{ji}}{\partial x_j} = -\frac{\partial \hat{V}_{ij}}{\partial x_i}$. Thus, it holds that

$$\begin{aligned} &\frac{\kappa_x}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} (v_j^T - v_i^T) \nabla_{x_j} \hat{V}_{ji} \\ &= \frac{\kappa_x}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} \left((v_j^T - v_{i_f}^T) \nabla_{x_j} \hat{V}_{ji} \right. \\ &\quad \left. + (v_i^T - v_{i_f}^T) \nabla_{x_i} \hat{V}_{ij} \right) \\ &= \kappa_x \sum_{j \in \mathcal{N}(i_f)} (v_j^T - v_{i_f}^T) \sum_{i \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} \nabla_{x_j} \hat{V}_{ji} \end{aligned} \quad (20)$$

Combining (19) and (20) yields that

$$\begin{aligned} \dot{H}(x_f, v_f) &= -\kappa_v \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f})^T (v_j - v_{i_f}) - \tilde{k}^T (C_{i_f}^F)^T C_{i_f}^F \tilde{k} \\ &\quad - \frac{\kappa_v}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{p \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} ((v_j - v_{i_f}) - (v_p - v_{i_f}))^T (v_j - v_p) \\ &= -\kappa_v \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f})^T (v_j - v_{i_f}) - \frac{\kappa_v}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{p \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} \\ &\quad (v_j - v_p)^T (v_j - v_p) - \tilde{k}^T (C_{i_f}^F)^T C_{i_f}^F \tilde{k} \end{aligned} \quad (21)$$

Therefore, $\dot{H}(x_f, v_f)$ is always nonpositive and $H(t) \leq H(0)$ for $t \geq 0$. From the definition of $H(t)$ in (18), it holds that $H(t) > \hat{V}_{ji}(t), \forall i, j \in \mathcal{V}_g$. Thus, $\hat{V}_{ji}(t) < H(0)$ for $t \geq 0$. According to the property 2) of \hat{V} , $\hat{V}_{ji}(t) > \bar{H}, \forall j \in \mathcal{N}(i)$ when $|x_{ji}| = 0$ and $|x_{ji}| = R$. Since the constant \bar{H} is chosen as in (16), it holds that $\bar{H} > H(0)$. Thus, $\hat{V}_{ji}(t) > \bar{H} > H(0)$ when $|x_{ji}| = 0$ and $|x_{ji}| = R$. This is contradict to $\hat{V}_{ji}(t) < H(0), \forall t \geq 0$. Hence, $|x_{ji}| \neq 0$ and $|x_{ji}(t)| \neq R, \forall t \geq 0$. This guarantees that the collision is avoided and no edge is lost between any two agents in \mathcal{V}_g .

Define the level set $\Omega_f \triangleq \{(x_f^T, v_f^T)^T \in \mathfrak{R}^{2(s+1) \times m} : H(t) \leq H(0), H(0) > 0\}$. By applying LaSalle's invariance principle, $(x_f^T, v_f^T)^T$ starting in Ω_f asymptotically converges to the largest invariant set inside the region $\mathcal{C} \triangleq \{(x_f^T, v_f^T)^T \in \Omega_f : \dot{H}(t) = 0\}$. According to (21), $\dot{H}(t) = 0$ holds if and only if $v_1 = v_2 = \dots = v_{i_f}$ and $C_{i_f}^F \tilde{k} = 0$. This implies that $\lim_{t \rightarrow \infty} (v_i - v_j) = 0, \forall i, j \in \mathcal{V}_g$. and $\lim_{t \rightarrow \infty} C_{i_f}^F \tilde{k} = 0$. Moreover, according to (12), $\lim_{t \rightarrow \infty} C_{i_f}^F = C_{i_f}^F/a$ where $a > 0$. Thus $\lim_{t \rightarrow \infty} C_{i_f}^F \tilde{k} = 0$.

In the following, we consider the error system $\dot{v}_j - \dot{v}_{i_f} = u_j - u_{i_f}$ for $j \in \mathcal{N}(i_f)$ at the point $v_1 = v_2 = \dots = v_{i_f}$. Obviously, it holds that $\dot{v}_j - \dot{v}_i = 0, \forall i, j \in \mathcal{V}_g$. Combining malicious controller (4) and controller (15), one has $\dot{v}_j - \dot{v}_{i_f} = -\kappa_v \sum_{j \in \mathcal{N}(i_f)} (v_{i_f} - v_j) - \sum_{p \in \mathcal{N}(j) \cap \mathcal{V}_g} (v_j - v_p) - \sum_{p \in \mathcal{N}(j) \cap \mathcal{V}_g} \nabla_{x_j} \hat{V}_{jp} + C_{i_f}^F \tilde{k}$. Note that $C_{i_f}^F \tilde{k} = 0$ at the point $v_1 = v_2 = \dots = v_{i_f} \forall j \in \mathcal{V}_g$. Thus, $\sum_{p \in \mathcal{N}(j) \cap \mathcal{V}_g} \nabla_{x_j} \hat{V}_{jp} = 0$. Define $\psi_{\mathcal{G}_f}(x_f) \triangleq (\dots, |x_{ij} - x_{ij}^*|, \dots)^T$ with $i, j \in \mathcal{V}_g$. Consider the error system in the compact form, one obtains that $R_{\mathcal{G}_f}^T(x_f) \xi(x_f) = 0$ where $\xi(x_f) \triangleq (\dots, \partial \hat{V}_{ij} / \partial |x_{ij} - x_{ij}^*|, \dots)$ and $R_{\mathcal{G}_f}^T(x_f) \triangleq \nabla_{x_f} \psi_{\mathcal{G}_f}(x_f)$ is the rigidity matrix. Since $\text{Rank}(R_{\mathcal{G}_f}(x_f)) = md - m(m+1)/2$ where m is the dimension and d is the vertex number of \mathcal{V}_g , it follows from the Rigidity Theory in [25] that $R_{\mathcal{G}_f}^T(x_f) \xi(x_f) = 0$ is equivalent to $\xi(x_f) = 0$. From the property 1) of \hat{V}_{ij} , we can deduce that $\partial \hat{V}_{ij} / \partial |x_{ij} - x_{ij}^*| = 0$ is equivalent to $|x_{ij} - x_{ij}^*| = 0, \forall i, j \in \mathcal{V}_g$. Hence, it holds that $x_{ij} \rightarrow x_{ij}^*$ as $t \rightarrow \infty$. Also, it yields from controllers (4) and (15) that $\dot{v}_{i_f} = u_{i_f} \rightarrow 0$ and $\dot{v}_j = u_j \rightarrow 0$ for $j \in \mathcal{N}(i_f)$.

Part B. As is shown in Fig. 2, all agents in Layer 2 can be viewed as the leaders of agents in Layer 3. Let $\bar{\mathcal{G}}$ be the direct graph

characterizing the information interaction among agents in \mathcal{V}_f and the transmission from agents in $\mathcal{V}_l(\mathcal{N}(i_f))$ to agents in \mathcal{V}_f . If there exists a directed path from agent $j \in \mathcal{N}(i_f)$ to agent $k \in \mathcal{V}_f$ in graph $\bar{\mathcal{G}}$, agent j is said to be a leader of agent k . Here, we prove that leader-follower flocking for agents in Layers 2-3 can be realized under controller (17) by analyzing the graph corresponding to each leader and all its followers. This together with the results in *Part A* yields that all the followers tend to the same velocity as that of the leaders.

Denote $\mathcal{L}(k)$ as the set of agent k 's leaders. Define the energy function $\Upsilon(x, v)$ as

$$\begin{aligned} \Upsilon(x, v) &\triangleq H(x_f, v_f) + \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{L}(i) \cap \mathcal{N}(i)} V_{ij} \\ &+ \frac{1}{2} \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} s_{\mathcal{L}(i)} V_{ip} \\ &+ \frac{1}{2} \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{L}(i)} (v_i - v_j)^T (v_i - v_j) \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} \frac{1}{2\gamma_{ij}} (\alpha_{ij} - \bar{\alpha})^2 \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \frac{s_{\mathcal{L}(i)}}{4\gamma_{ip}} (\alpha_{ip} - \bar{\alpha})^2 \end{aligned} \quad (22)$$

where $s_{\mathcal{L}(i)}$ denotes the number of agents in set $\mathcal{L}(i)$. Constant $\bar{\alpha}$ will be designed later.

Note that the graph of agents in Layer 3 is undirected, thus $s_{\mathcal{L}(i)} = s_{\mathcal{L}(j)}$, $\forall i \in \mathcal{V}_f, j \in \mathcal{N}(i) \cap \mathcal{V}_f$. Therefore, the derivative of Υ is

$$\begin{aligned} \dot{\Upsilon} &= \dot{H} + \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{L}(i) \cap \mathcal{N}(i)} \dot{x}_{ij}^T \nabla_{x_{ij}} V_{ij} + \frac{1}{2} \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} s_{\mathcal{L}(i)} \dot{x}_{ip}^T \nabla_{x_{ip}} V_{ip} \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{L}(i)} (v_i - v_j)^T (\dot{v}_i - \dot{v}_j) \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} (\alpha_{ij} - \bar{\alpha}) |v_i - v_j|_1 \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \frac{s_{\mathcal{L}(i)}}{2} (\alpha_{ip} - \bar{\alpha}) |v_i - v_p|_1 \end{aligned}$$

Applying controller (17), one obtains that

$$\begin{aligned} \dot{\Upsilon} &= \dot{H}(x_f, v_f) + \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{L}(i) \cap \mathcal{N}(i)} (v_i - v_j)^T \nabla_{x_{ij}} V_{ij} \\ &+ \frac{1}{2} \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} s_{\mathcal{L}(i)} \dot{x}_{ip}^T \nabla_{x_{ip}} V_{ip} + \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{L}(i)} (v_i - v_j)^T \\ &\times \left(- \sum_{j \in \mathcal{L}(i) \cap \mathcal{N}(i)} \alpha_{ij} \text{sgn}(v_i - v_j) - \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \alpha_{ip} \text{sgn}(v_i - v_p) \right. \\ &\left. - \sum_{j \in \mathcal{L}(i) \cap \mathcal{N}(i)} \nabla_{x_{ij}} V_{ij} - \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \nabla_{x_{ip}} V_{ip} - u_j \right) \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} (\alpha_{ij} - \bar{\alpha}) |v_i - v_j|_1 \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \frac{s_{\mathcal{L}(i)}}{2} (\alpha_{ip} - \bar{\alpha}) |v_i - v_p|_1 \end{aligned}$$

For convenience, label the agents in $\mathcal{N}(i_f)$ who have neighbors in \mathcal{V}_f as 1 to ω . If there exists a directed path from agent $j, j \in \{1, \dots, \omega\}$ to some agents in \mathcal{V}_f , denote the set of these agents as $F(j)$. Note that $F(j) \subseteq \mathcal{V}_f$, and the leaders of k, p are same if $k, p \in \mathcal{V}_f$ are

neighbors. Therefore, it yields that

$$\begin{aligned} \dot{\Upsilon} &= \dot{H} + \frac{1}{2} \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} s_{\mathcal{L}(i)} \dot{x}_{ip}^T \nabla_{x_{ip}} V_{ip} \\ &- \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} \alpha_{ij} |v_i - v_j|_1 - \sum_{j=1}^{\omega} \sum_{i \in F(j)} (v_i - v_j)^T \\ &\times \sum_{p \in F(j) \cap \mathcal{N}(i)} \alpha_{ip} \text{sgn}((v_i - v_j) - (v_p - v_j)) \\ &- \sum_{j=1}^{\omega} \sum_{i \in F(j)} (v_i - v_j)^T \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \nabla_{x_{ip}} V_{ip} \\ &- \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} (v_i - v_j)^T u_j - \sum_{j=1}^{\omega} \sum_{i \in F(j)} (v_i - v_j)^T u_j \\ &+ \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} (\alpha_{ij} - \bar{\alpha}) |v_i - v_j|_1 \\ &+ \sum_{j=1}^{\omega} \sum_{i \in F(j)} (v_i - v_j) \sum_{p \in F(j) \cap \mathcal{N}(i)} (\alpha_{ip} - \bar{\alpha}) \\ &\times \text{sgn}((v_i - v_j) - (v_p - v_j)) \end{aligned} \quad (23)$$

Since V_{ip} is symmetric with respect to x_{ip} and $x_{ip} = (x_i - \chi) - (x_p - \chi)$, $\forall \chi \in \mathbb{R}^m$, it holds that

$$\begin{aligned} &\frac{1}{2} \sum_{i \in \mathcal{V}_f} \sum_{p \in \mathcal{V}_f \cap \mathcal{N}(i)} \dot{x}_{ip}^T \nabla_{x_{ip}} V_{ip}(x_{ip}) \\ &= \frac{1}{2} \sum_{j=1}^{\omega} \sum_{i \in F(j)} \sum_{p \in F(j) \cap \mathcal{N}(i)} ((v_i - v_j) - (v_p - v_j)) \\ &\quad \times \nabla_{x_{ji}} \hat{V}_{ip}((x_i - x_j) - (x_p - x_j)) \\ &= \sum_{j=1}^{\omega} \sum_{i \in F(j)} (v_i - v_j) \sum_{p \in F(j) \cap \mathcal{N}(i)} \nabla_{x_{ip}} \hat{V}_i(x_{ip}) \end{aligned} \quad (24)$$

Since $u_j(t)$ is continuous for $t \in [0, \infty)$ and it is proved in *Part A* that $\lim_{t \rightarrow \infty} u_j = 0$, it holds that $u_j(t)$ is bounded for $t \in [0, \infty)$. Denote the bound as μ such that $|u_j(t)|_1 \leq \mu, \forall j \in \mathcal{N}(i_f)$. Substituting (24) into (23) yields

$$\begin{aligned} \dot{\Upsilon} &\leq \dot{H} + \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} |u_j| |v_i - v_j| \\ &+ \sum_{j=1}^{\omega} \sum_{i \in F(j)} |u_j| |v_i - v_j| - \sum_{i \in \mathcal{V}_f} \sum_{j \in \mathcal{N}(i_f) \cap \mathcal{N}(i)} \bar{\alpha} |v_i - v_j|_1 \\ &- \sum_{j=1}^{\omega} \sum_{i \in F(j)} (v_i - v_j)^T \sum_{p \in F(j) \cap \mathcal{N}(i)} \bar{\alpha} \text{sgn}((v_i - v_j) - (v_p - v_j)) \end{aligned}$$

Define the number of agents in $F(j), j \in \{1, \dots, \omega\}$ as $s_{F(j)}$. Define \tilde{v}_j as a column stack vector of $(v_i - v_j), i \in F(j)$. Let \mathcal{G}_j be the undirected graph characterizing the interaction among the $s_{F(j)}$ followers of leader j with the associated Laplacian matrix $L_j \triangleq D_j D_j^T$. Note that by definition, L_j is symmetric positive semi-definite. Let $\bar{\mathcal{G}}_j$ be the directed graph characterizing the interaction among leader j and its followers. Let the edge weight $a_{ij} = 1$ if leader j is a neighbor of follower i and $a_{ij} = 0$ otherwise. Define $\Lambda_j \triangleq \text{diag}(a_{i_1 j}, \dots, a_{i_{s_{F(j)}} j}), i_k \in F(j), k = 1, \dots, s_{F(j)}$. Note that $\Lambda_j^2 = \Lambda_j$ because $a_{ij}, i \in F(j)$ is either 1 or 0. Therefore, it

holds

$$\begin{aligned} \dot{\Upsilon} &\leq \dot{H} + \sum_{j=1}^{\omega} \mathbf{1}_{s_{F(j)}} \otimes \mu |\check{v}_j| - \sum_{j=1}^{\omega} \bar{\alpha} |\Lambda_j \otimes I_m \check{v}_j| \\ &\quad - \sum_{j=1}^{\omega} \bar{\alpha} |D_j^T \otimes I_m \check{v}_j| \end{aligned}$$

Define the leader-follower topology matrix associated with graph $\bar{\mathcal{G}}_j$ as $R_j \triangleq L_j + \Lambda_j$, $j \in \{1, \dots, \omega\}$. According to Lemma 2, R_j is symmetric positive definite. Based on (21) and the fact that $|\cdot| \leq |\cdot|_1$ for any vector, one obtains

$$\begin{aligned} \dot{\Upsilon} &\leq -\kappa_v \sum_{j \in \mathcal{N}(i_f)} (v_j - v_{i_f})^T (v_j - v_{i_f}) \\ &\quad - \frac{\kappa_v}{2} \sum_{j \in \mathcal{N}(i_f)} \sum_{p \in \mathcal{N}(i_f) \cap \mathcal{N}(j)} (v_j - v_p)^T (v_j - v_p) \\ &\quad - \tilde{k}^T (C_{i_f}^F)^T C_{i_f}^F \tilde{k} - \sum_{j=1}^{\omega} (\bar{\alpha} \sqrt{\lambda_{\min}(R_j)} - \bar{\mu}) |\check{v}_j| \quad (25) \end{aligned}$$

where $\bar{\mu} \triangleq \max_{k \in \{s_{F(1)}, s_{F(2)}, \dots, s_{F(\omega)}\}} \{\mathbf{1}_k \otimes \mu\}$. If $R_j(t)$, changes at some time, there exists $t_1 > 0$ such that $R_j(t) = R_j(0)$, $\forall t \in [0, t_1)$. By designing $\bar{\alpha} > \bar{\mu} / \min_{j \in \{1, \dots, \omega\}} \{\sqrt{\lambda_{\min}(R(j)(0))}\}$, one has $\dot{\Upsilon}(t) \leq 0$ for $t \in [0, t_1)$. Since $V_{ik}(t)$, $i \in \mathcal{V}'$, $k \in \mathcal{N}(j) \cap \mathcal{V}'$ is continuous, we can conclude that $V_{ik}(t_1) \leq \Upsilon(t_1)$. From the definition of V_{ik} in (2), it follows that there is no collision and no edge in the graph $\bar{\mathcal{G}}_j(0)$ will be lost for $t \in [0, t_1]$. Therefore, the only possibility that $R_j(t)$ changes at $t = t_1$ is that some edges are added in the graph, which means that $\bar{\mathcal{G}}_j(0)$ is a subgraph of $\bar{\mathcal{G}}_j(t_1)$. Then it yields from Lemma 2.2 in [24] that $R_j(0) \leq R_j(t_1)$, $\forall j \in \{1, \dots, \omega\}$ and thus $\lambda_{\min}(R_j(0)) \leq \lambda_{\min}(R_j(t_1))$. Following the same argument, if $R_j(t)$ changes at $t = t_i > t_1$, $i = 1, \dots$, one can have the same conclusion. Therefore, it holds that $\bar{\mu} / \min_{j \in \{1, \dots, \omega\}} \{\sqrt{\lambda_{\min}(R(j)(0))}\} \geq \bar{\mu} / \min_{j \in \{1, \dots, \omega\}} \{\sqrt{\lambda_{\min}(R(j)(t))}\}$ for all $t \geq 0$. And there is no collision and also no edge in the graph $\bar{\mathcal{G}}_j(0)$ is lost for all $t \geq 0$. Thus, $\dot{\Upsilon}(t) \leq 0$, $\forall t \geq 0$.

Combining (25) with the analysis in Part A, it holds that $\lim_{t \rightarrow \infty} (v_j - v_{i_f}) = 0$ and $\lim_{t \rightarrow \infty} (v_i - v_j) = 0$, $\forall i \in F(j)$, $j \in \{1, \dots, \omega\}$. Assumption 2 indicates that there exists at least one leader in $\{1, \dots, \omega\}$ for each agent in \mathcal{V}_f when $t = 0$. Since no edge in $\bar{\mathcal{G}}_j$, $\forall j \in \{1, \dots, \omega\}$ is lost for $t \geq 0$, all agents in \mathcal{V}_f are followers of agents in $\{1, \dots, \omega\}$ for $t \geq 0$. This further leads to $\lim_{t \rightarrow \infty} (v_i - v_j) = 0$, $\forall i, j \in \mathcal{V}$. This completes the proof. \square

Remark 4 : According to the proof of Theorem 1, $\lim_{t \rightarrow \infty} (x_{ij} - x_{i_j}^*) = 0$, $\forall i, j \in \mathcal{V}_g$ holds. Since $|x_{i_f j}^*| = \bar{\delta}_{i_f j} < R/2$, $\forall j \in \mathcal{N}(i_f)$, it holds $|x_{pk}^*| = |x_{pi_f}^* + x_{i_f k}^*| \leq |x_{i_f p}^*| + |x_{i_f k}^*| < R$, $\forall p, k \in \mathcal{N}(i_f)$. Also, no edge in \mathcal{V}_g is lost under the controller (15). This together with Assumption 4 that any two agents in $\mathcal{N}(i_f)$ are neighbors at the initial time guarantees that any two agents in $\mathcal{N}(i_f)$ are always neighbors for all $t \geq 0$. Therefore, the local information exchange among all the neighbors of the malicious agent can be obtained as they can interact with each other. If Assumption 4 is not satisfied, this information can be achieved in virtue of a local communication network among agents in Layer 2 [22]. Such a local network can be built conveniently if it does not exist, since these agents are close to each other. With this local network, $\bar{\delta}_{i_f j}$ can be any desired value in $(0, R)$. \square

Remark 5 : The main idea of the geometric configuration control (15) is to contain the malicious agent by ‘‘pulling’’ its neighbors to the desired geometric shape. In the controller, the first term is to urge the agents to reach the same common velocity. The second term is

to let the agents approach to the desired configuration to contain the malicious agent. The last term is to compensate for the influence of the malicious agent reacting on its neighbors. \square

IV. EXPERIMENTAL RESULT

In this section, the experimental result is presented to illustrate the effectiveness of the proposed flocking control scheme in the above section.

A semi-physical experimental platform of an Unmanned Aerial Vehicle(UAV) swarm has been set up based on 40 Raspberry Pi computers. The dynamics and controller of each UAV are simulated by 2 Raspberry Pi computers, respectively. Specifically, the flight control system model of UAV in the platform and the transformation method between the UAV model and model (1) are from Ref. [21]. Fig. 3 is the picture of the UAV swarm semi-physical platform, which consists of 4 parts: Raspberry Pi computers, a thrust lever, a data analysis terminal and a flight display terminal.

In the experiment, we consider a 2-dimensional swarm of 13 UAVs (UAVs 0-12), including a malicious agent (UAV 6) under controller (4)-(5) with $k_a = 0$, $k_r = 450000$ and $k_v = 0.8$. Define the velocity of UAV $i \in \{0, 1, \dots, 12\}$ as $v_i \triangleq (v_{xi}, v_{yi})$ where v_{xi} and v_{yi} are velocities in x-dimension and y-dimension, respectively. The control inputs of UAV i is the banking angle Φ_i , lift L_i and engine thrust T_i . The initial ground velocity V_i , $i \in \{0, 1, \dots, 12\}$ of UAV i is taken randomly from $(27, 35)m/s$ and heading angle χ_i is taken from $(\pi/6, \pi/4)$ rad. The initial flight path angle is 0. According to the model transformation in [21], $v_{xi} \triangleq V_i \cos(\gamma_i) \cos \chi_i$ and $v_{yi} \triangleq V_i \cos(\gamma_i) \sin \chi_i$. Let the communication distance be $R = 18\sqrt{2}m$, thus $\delta = 18m$. Let the desired distance between the malicious agent and its neighbors be $\bar{\delta} = 12m < R/2$. Apply controller (15) with $\kappa_v = 6$ and $\kappa_x = 2$ to UAVs 2, 5, 7 and 10. Apply controller (17) with $\gamma_{kp} = 1$ to UAVs 0, 1, 3, 4, 8, 9, 11 and 12. The experimental result presented in Figs. 4-5 shows that all UAVs tend to a common velocity and all the control efforts tend to 0. The malicious UAV 6 is contained, and the distances between it and its neighbors tend to 12m as expected and the configuration tends to the desired one as shown in Fig. 6.

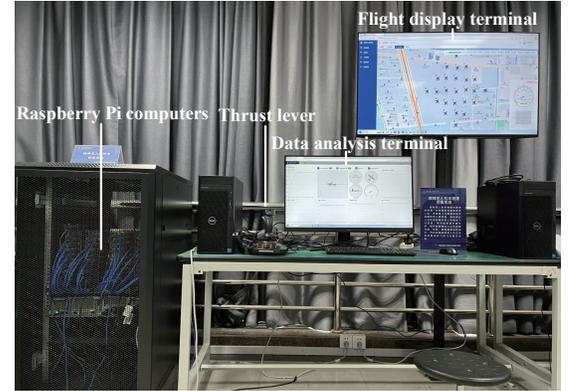


Fig. 3. Overview of UAV swarm semi-physical experimental platform.

V. CONCLUSION

This paper, for the first time, considers the flocking control with a malicious agent, and the proposed hierarchical geometric configuration based flocking control method is applied to a swarm with a malicious agent. The new result enriches the conventional flocking control theory. In the future, by combining the switching system theory and the proposed parameter estimation framework, the malicious agent with changeable parameters will be taken into

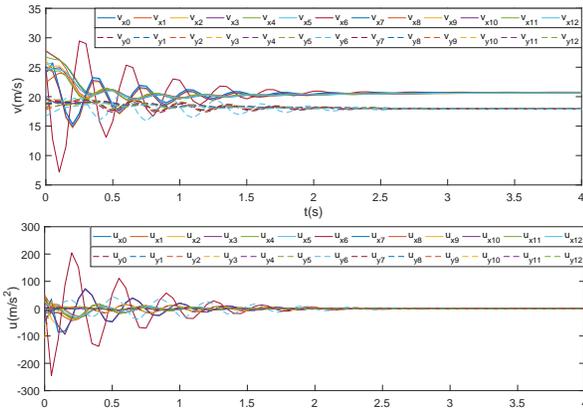


Fig. 4. Trajectories of velocities and control efforts of UAVs 0-12.

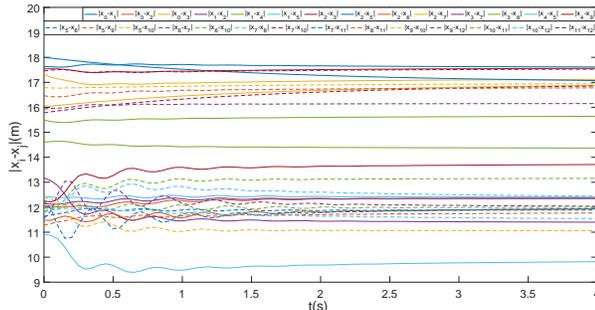


Fig. 5. Trajectories of distances between the neighboring UAVs.

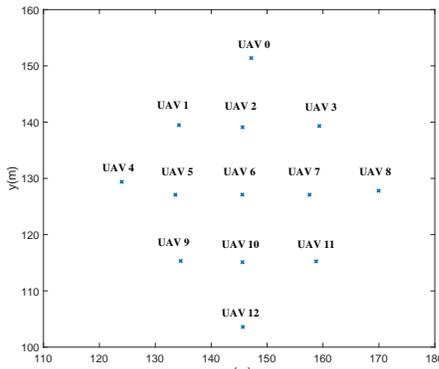


Fig. 6. Flocking patterns at $t = 4s$ where x and y present positions in x and y dimensions, respectively.

consideration. Moreover, this new result will be extended to more cases: one malicious agent acts selectively on a part of its neighbors, or multiple malicious agents existing in the swarm. Further studies will also focus on the application of the containment method to multi-agent with nonlinear or other complex dynamics.

REFERENCES

- [1] C. W. Reynolds. "Flocks, herds and schools: A distributed behavioral model," in *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques*, 1987, pp. 25-34.
- [2] R. Olfati-Saber. "Flocking for multi-agent dynamic systems: Algorithms and theory," *IEEE Transactions on Automatic Control*, vol. 51, no. 3, pp. 401-420, Mar. 2006.
- [3] H. G. Tanner, A. Jadbabaie, G. J. Pappas. "Flocking in fixed and switching networks," *IEEE Transactions on Automatic Control*, vol. 52, no. 5, pp. 863-868, May 2007.
- [4] V. Gazi and K. M. Passino. *Swarm stability and optimization*. Springer Science & Business Media, 2011.

- [5] G. H. Wen, Z. S. Duan, H. S. Su, G. R. Chen, W. W. Yu. "A connectivity-preserving flocking algorithm for multi-agent dynamical systems with bounded potential function," *IET Control Theory & Applications*, vol. 6, pp. 813-821, Jan. 2011.
- [6] Y. Dong, J. Huang. "Flocking with connectivity preservation of multiple double integrator systems subject to external disturbances by a distributed control law," *Automatica*, vol. 55, pp. 197-203, May 2015.
- [7] G. Fedele, L. D'Alfonso, and A. Bono. "A discrete-time model for swarm formation with coordinates coupling matrix," *IEEE Control Systems Letters*, vol. 4, no. 4, pp. 1012-1017, May 2020.
- [8] A. Bono, G. Fedele, and G. Franzè. "A swarm-based distributed model predictive control scheme for autonomous vehicle formations in uncertain environments," *IEEE Transactions on Cybernetics*, pp. 1-11, May 2021.
- [9] H. Yang, Q. L. Han, X. H. Ge, L. Ding, Y. H. Xu, B. Jiang, D. H. Zhou. "Fault-tolerant cooperative control of multiagent systems: A survey of trends and methodologies," *IEEE Transactions on Industrial Informatics*, vol. 16, pp. 4-17, Jan. 2020.
- [10] D. Ding, Q. L. Han, Y. Xiang, X. Ge, X. M. Zhang. "A survey on security control and attack detection for industrial cyber-physical systems," *Neurocomputing*, vol. 275, pp. 1674-1683, Jan. 2018.
- [11] M. S. Sanders, E. J. McCormick. "Human factors in engineering and design," *Industrial Robot*, vol. 25, pp. 153, 1998.
- [12] S. Dhama. *The foundations of behavioral economic analysis*. New York, USA: Oxford University Press, 2016.
- [13] S. Yazdani, M. Haeri. "Robust adaptive fault-tolerant control for leader-follower flocking of uncertain multi-agent systems with actuator failure," *ISA transactions*, vol. 71, pp. 227-234, Nov. 2017.
- [14] Z. Feng, G. Hu. "Connectivity-preserving flocking for networked Lagrange systems with time-varying actuator faults," *Automatica*, vol. 109, pp. 108509, Nov. 2019.
- [15] K. L. Hobbs, C. Cargal, E. Feron, R. S. Burns. "Early safety analysis of manned-unmanned team system," *2018 AIAA Information Systems-AIAA Infotech @ Aerospace*, pp. 1-15, Jan. 2018.
- [16] L. An, G. H. Yang. "Byzantine-resilient distributed state estimation: A min-switching approach," *Automatica*, vol. 129, pp. 109664, Jul. 2021.
- [17] K. Saulnier, D. Saldaña, A. Prorok, G. J. Pappas, and V. Kumar. "Resilient flocking for mobile robot teams," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 1039-1046, Apr. 2017.
- [18] S. Gil, C. Baykal, and D. Rus. "Resilient multi-agent consensus using wi-fi signals," *IEEE Control Systems Letters*, vol. 3, no. 1, pp. 126-131, Jan. 2019.
- [19] F. Mallmann-Trenn, M. Cavorsi, and S. Gil. "Crowd vetting: Rejecting adversaries via collaboration with application to multirobot flocking," *IEEE Transactions on Robotics*, pp. 1-20, Feb. 2021.
- [20] Y. Shang. "Consensus of hybrid multi-agent systems with malicious nodes," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 67, no. 4, pp. 685-689, May 2019.
- [21] P. K. Menon, G. D. Sweriduk, B. Sridhar. "Optimal strategies for free-flight air traffic conflict resolution," *Journal of Guidance, Control, and Dynamics*, vol. 22, no. 2, pp. 202-211, Mar. 1999.
- [22] H. Yang, B. Jiang, M. Staroswiecki, Y. M. Zhang. "Fault recoverability and fault tolerant control for a class of interconnected nonlinear systems," *Automatica*, vol. 54, pp. 49-55, Apr. 2015.
- [23] J. Mei, W. Ren, G. Ma. "Distributed coordinated tracking for multiple euler-lagrange systems," in *49th IEEE Conference on Decision and Control*, 2010, pp. 3208-3213.
- [24] S. Ghapani, J. Mei, W. Ren, Y. Song. "Fully distributed flocking with a moving leader for lagrange networks with parametric uncertainties," *Automatica*, vol. 67, pp. 67-76, May 2016.
- [25] B. Hendrickson. "Conditions for unique graph realizations," *SIAM Journal on Computing*, vol. 21, no. 1, pp. 65-84, Jul. 1992.