

# A Cyber-Physical Game Framework for Secure and Resilient Multi-Agent Autonomous Systems

Zhiheng Xu and Quanyan Zhu

Department of Electrical and Computer Engineering, Polytechnic School of Engineering,  
New York University, Brooklyn, 11201, USA. E-mail: {zx383, qz494}@nyu.edu.

**Abstract**—The increasing integration of autonomous systems with publicly available networks exposes them to cyber attackers. An adversary can launch a man-in-the-middle attack to gain control of the system and inflict maximum damages with collision and suicidal attacks. To address this issue, this work establishes an integrative game and control framework to incorporate security into the automatic designs, and take into account the cyber-physical nature and the real-time requirements of the system. We establish a cyber-physical signaling game to develop an impact-aware cyber defense mechanism and leverage model-predictive control methods to design cyber-aware control strategies. The integrative framework enables the co-design of cyber-physical systems to minimize the inflicted systems, leading to online updating the cyber defense and physical layer control decisions. We use unmanned aerial vehicles (UAVs) to illustrate the algorithm, and corroborate the analytical results in two case studies.

## I. INTRODUCTION

Autonomous System (AS) technologies for civil applications, e.g., Amazon drones for package delivery [1], Google self-driving cars [2], are increasingly integrated with the existing publicly available networks, such as the Internet, with temporary ad-hoc wireless and satellite networks to send video, audio and other sensor or actuation data to remote operators [3], [4]. This exposes AS to cyber attacks due to the existing vulnerabilities of the cyber infrastructure, and the lack of security considerations in AS design. For example, an attacker can hijack the communication channel between the control station and the AS, and send wrong reference signals or control commands to turn the AS into a weapon, and command it to smash into a building or collide with other ASs.

To address this issue, several features of the ASs have to be taken into account. The first and foremost is the cyber-physical nature of the system. The performance of physical layer control system is significantly dependent on the integrity and availability of the sensor, actuator, and reference signal data, which can be communicated remotely from the control station. The success of a cyber attack will not only compromise data security, but also result in instability of the control system, and consequently unanticipated man-made disasters. The defense mechanism has to be designed to be aware of the impact of the attack on both cyber and physical layer of the system. The second feature is the real-time system requirement. ASs are

constrained by their computational and communication resources. Traditional cryptographic solutions can be expensive and ineffective for sophisticated attacks, such as Stuxnet [5], Duqu [6] and other advanced persistent attacks (APTs) [7]. Hence defense mechanisms should be light-weight, and create minimum delay and computational overhead so that the real-time system performance requirements are met.

In this paper, we aim to design security mechanisms for two adversary models. One is the suicidal attack and the other is the collision attack. An attacker can launch a spoofing attack [8] or man-in-the-middle attack [9] to modify the input reference signal communicated between the control station and the AS. In the suicidal attack, the adversary aims to collide the AS with an obstacle such as a building or a bridge. In the collision attack, the adversary aims to make a group of ASs collide into each other. These attacks will become a lethal weapon when a swarm of ASs can be taken over by an adversary. Hence it is particularly important to address the two attack models in the context of multi-agent ASs.

To this end, we develop a cyber-physical signaling games to capture the information asymmetry and the multi-stage behaviors between the attacker and the multiple ASs, and establish a model predictive control (MPC) based control system model that connects the cyber and physical layers of the ASs. This holistic approach to AS security is essential to address the two main features of the ASs, enabling the co-design of impact-aware cyber defense mechanisms, and the cyber-aware control strategies. The proposed mechanism integrates the cyber and physical layers that couple the updates of the control and defense decisions.

To verify our mechanism, we use unmanned aerial vehicles (UAVs) to corroborate the analytical results in two case studies. Our results show that the cyber-physical game admits a Perfect Bayesian Nash Equilibrium (PBNE), and the equilibrium strategies lead to the protection of ASs from collisions. One significant advantage is that this PBNE can be achieved quickly using the proposed mechanism, deterring an adversary to launch successful attacks. The main contributions of this paper are summarized as follows:

- We develop a cyber-physical game that connects the minimax game at the physical control layer and the signaling game at the cyber layer of ASs. The integrated game captures the in between cyber and physical layers in the context of cyber threats.
- We identify the cyber-physical attack models of multi-

This work was in part supported by the NSF grants (EFMA-1441140, SES-1541164) and a grant from NYU Research Challenge Fund.

agent ASs that can inflict a significant and lethal impact on the society, and provide control and game-theoretic solutions to reduce their impact.

- We investigate one-sender and multi-receiver signaling game and define a PBNE of the cyber-physical signaling game. We characterize the separating and pooling equilibrium strategies in closed form.

The paper is organized as follows. Section 2 presents the problem statement. In Section 3, we analyze the cyber-physical signaling game and develop a theorem for our mechanism. Simulation results are presented in Section 4. Finally, Section 5 concludes this paper.

## II. PROBLEM STATEMENT

In this section, we consider a cyber-physical architecture for autonomous systems (ASs) and describe the system dynamics of ASs using a discrete-time linear system model. Based on the model, a min-max model predictive control (MPC) problem is formulated to handle the worst-case disturbances. We identify a class of attack models to describe the impact of an adversary on the system, and present two critical scenarios for investigation. To enhance the resilience to the attacks, each AS is equipped with a local sensor that can sense objects within a limited range. Finally, we introduce a signaling game model to capture the behaviors between ASs and an adversary, and characterize its perfect Bayesian Nash equilibrium (PBNE) solutions.

### A. The dynamic Model and the MPC problem

Consider a multi-agent AS with  $n$  agents and a control station (CS) that sends reference signals to each agent through wireless networks. The control objective of the agent  $i$  is to track a given reference signal  $r^i$  from the CS.

At each time instant  $k \in \mathbb{Z}_+$ , we use a linear discrete-time state-space model to describe the dynamics of the  $i^{\text{th}}$  AS, which is given as follows:

$$\begin{aligned} x_{k+1}^i &= A^i x_k^i + B^i u_k^i + w_k^i, \\ z_k^i &= H^i x_k^i, \\ x^i(0) &= x_0^i, \quad i \in \mathcal{N}, \end{aligned} \quad (1)$$

where  $\mathcal{N} \triangleq \{1, 2, \dots, n\}$  is the index set of the agents;  $x_k^i \in \mathbb{R}^{n_x}$  is state vector of size  $n_x$ ;  $u_k^i \in \mathbb{R}^{n_u}$  is the control vector of size  $n_u$ ;  $z_k^i \in \mathbb{R}^{n_z}$  is a controlled output vector of size  $n_z$ ;  $w_k^i \in \mathbb{R}^{n_x}$  is a unknown disturbance vector;  $x_0^i \in \mathbb{R}^{n_x}$  is a given initial condition;  $A^i$ ,  $B^i$ ,  $C^i$ , and  $H^i$  are constant matrices with appropriate dimensions for AS  $i$ . The system is subject to state and the control constraints:

$$u^i \in \mathcal{U}, \text{ and } x^i \in \mathcal{X},$$

where  $\mathcal{X} \subset \mathbb{R}^{n_x}$  and  $\mathcal{U} \subset \mathbb{R}^{n_u}$  are compact sets containing the origin. We assume that the unknown disturbance  $w_k^i$  belongs to a compact set  $\mathcal{W} \subset \mathbb{R}^{n_x}$ . Throughout this paper,  $\|\cdot\|$  denotes the Euclidean norm.

In this paper, we assume that the controlled output  $z_k^i$  is the position of the agent  $i$  at time  $k$ . Each AS aims to track its reference trajectory  $r_k^i$  given by the CS, where  $r_k^i \in \mathcal{R}_k^i \subset$

$\mathbb{R}^{n_z}$ , and  $\mathcal{R}_k^i$  is a feasible set for  $r_k^i$ . To achieve its trajectory-tracking goal, each AS solves a min-max MPC problem to compute its control inputs. MPC yields a moving-horizon strategy, and control inputs are computed at every sampling time. The periodical property of MPC makes it tractable to handle constraints on control inputs or system states [10].

Based on the system model (1), the min-max MPC problem  $\mathcal{P}_k^i$  of AS  $i$  at time  $k$  is formulated as

$$\begin{aligned} \mathcal{P}_k^i : \quad & \min_{\hat{u}_k^i \in \mathcal{U}^N} \max_{\hat{w}_k^i \in \mathcal{W}^N} J_c(x_k^i, r_k^i, \hat{u}_k^i, \hat{w}_k^i) \\ & \triangleq \sum_{\tau=0}^{N-1} h(\hat{x}_{k+\tau|k}^i, \hat{u}_{k+\tau|k}^i, r_k^i), \end{aligned}$$

subject to

$$\begin{aligned} \hat{x}_{k+\tau|k}^i &= A^i \hat{x}_{k+\tau-1|k}^i + B^i \hat{u}_{k+\tau-1|k}^i + \hat{w}_{k+\tau-1|k}^i, \\ \hat{u}_k^i &= (\hat{u}_{k|k}^i, \dots, \hat{u}_{k+N-1|k}^i), \\ \hat{w}_k^i &= (\hat{w}_{k|k}^i, \dots, \hat{w}_{k+N-1|k}^i), \\ \hat{x}_{k|k}^i &= x_k^i, \\ \hat{x}_{k+\tau|k}^i &\in \mathcal{X}, \quad \hat{u}_{k+\tau|k}^i \in \mathcal{U}, \quad \hat{w}_{k+\tau|k}^i \in \mathcal{W}, \\ &\forall \tau = 0, \dots, N-1, \end{aligned}$$

where  $h : \mathcal{X} \times \mathcal{U} \times \mathcal{R}^i \rightarrow \mathbb{R}$  is the stage cost function defined as  $h(\hat{x}^i, \hat{u}^i, r^i) = \|H^i \hat{x}^i - r^i\|^2 + \gamma \|\hat{u}^i\|^2$ , and  $\gamma$  is a tuning parameter;  $N$  is the horizon-window length;  $\hat{x}_{k+\tau|k}^i$  is the estimate value of  $x_{k+\tau}^i$  given the feedback state  $x_k^i$ . At each time instant  $k$ , the MPC problem  $\mathcal{P}_k^i$  is solved repeatedly, but only the first control  $\hat{u}_{k|k}^i$  is applied to the system [11].

### B. Attack Models

Many attack models on networked control systems have been discussed in the literature [12]. Potential threats on ASs can come from different layers of the system including physical-layer jamming [13] and spoofing [8], falsified data injection [14], and node capturing attacks [15]. This paper focuses on *Man-In-The-Middle* (MITM) type of attacks between the control station and multiple ASs. In the MITM attack [16], the adversary first blocks valid messages, then creates independent connections with both entities. By doing so, the adversary gains control, while making the remote operators and robots believe they are talking to each other. These attacks are dangerous since the adversary becomes a part of the process and exploits system properties to maximize the damage. Three common types of MITM attack are replay attack, message spoofing attack, and message dropping attack. In replay attack [17], an adversary replays the messages that he intercepts between a sender and a receiver. In the message spoofing attack [18], an adversary modifies the message that is sent to the AS. In the message dropping attack [19], the attacker can delay or drop the message from the operator to the AS. These attack schemes can lead to the following scenarios as a consequence:

- **Suicidal Attack:** The adversary aims to crash a group of ASs to a static obstacle, e.g. a building, power substation and bridge.

- **Collision Attack:** The adversary mounts an attack to misinform a group of moving ASs to make them collide into each other.

These two scenarios, if successful, not only create serious damages to ASs but also threats to the existing critical infrastructures. In addition, when the number of AS is large, the two attacks can turn ASs into a lethal weapon against our nation and creates unprecedented impact on our society.

In this paper, we assume that an adversary can fabricate a fake reference signal  $r^i$  to deviate agent  $i$  from its real trajectory  $\bar{r}^i$  to achieve Suicidal Attack (SA) or Collision Attack (CA). To withstand these attacks, we will develop a mechanism to determine whether the agent  $i$  should accept a given reference  $r^i$ .

### C. Cyber Layer Signaling Game Model

Game theory deals with strategic interactions among multiple decision makers. It has been widely applied in the field of cyber security [20]. In this paper, we integrate a game-theoretic approach with the MPC scheme to tackle the cyber-physical security problem.

Two types of players exist in our scenario: One is the CS (or sender), denoted as  $S$ , and the others are ASs (or receivers), denoted as  $R^i$ ,  $i \in \mathcal{N}$ . The CS and all the ASs constitute the players of the game. The CS sends a message to each AS, and AS chooses an action based on the message. The AS does not know whether the message is sent from a normal or malicious sender. Due to these unique characteristics, we use a signaling game method to capture the information asymmetry and multi-stage behaviors of these players.

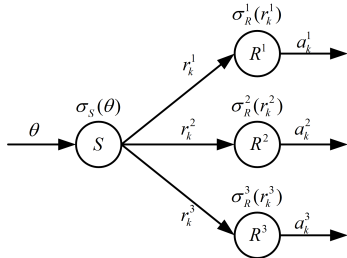


Fig. 1. A example of the signaling game model with one sender and three receivers: The sender, who has a private information  $\theta$ , sends a message  $\mathbf{r}_k = \{r_k^1, r_k^2, r_k^3\}$  using its strategy  $\sigma_S(\theta)$ ; the receiver  $R^i$  observes  $r_k^i$ , and chooses an action  $a_k^i$  using its strategy  $\sigma_R^i(r_k^i)$ .

$S$  has a private information called type  $\theta \in \Theta$ , where  $\Theta = \{\theta_0, \theta_1\}$  is the type space of  $S$ .  $\theta_0$  denotes that  $S$  is normal, while  $\theta_1$  denotes that  $S$  is malicious. Knowing its type  $\theta$ ,  $S$  sends a group of messages  $\mathbf{r} = (r^1, \dots, r^n)$  to each  $R^i$ , where  $r^i$  belongs to a compact set  $\mathcal{R}^i \subset \mathbb{R}^{n_z}$ . After observing  $r^i$ ,  $R^i$  chooses an action  $a^i \in \mathcal{A}$ , where  $\mathcal{A} = \{0, 1\}$  is the action set.  $a^i = 1$  means that  $R^i$  trusts  $S$  and accepts  $r^i$  as its reference, whereas  $a^i = 0$  means that  $R^i$  does not trust  $S$  and rejects  $r^i$ ; Sender has the cost function  $c_S : \mathcal{R} \times \mathcal{A} \times \Theta \rightarrow \mathbb{R}$  where  $\mathcal{R} = \prod_{i \in \mathcal{N}} \mathcal{R}^i$ . The receiver  $R^i$  has the cost function  $c_R^i : \mathcal{R}^i \times \mathcal{A} \times \Theta \rightarrow \mathbb{R}$ . Fig. 1 illustrates an example of signaling game with a sender and three receivers.

$R^i$  has a prior belief  $p^i(\theta)$  about the type  $\theta$  of  $S$ . Sender  $S$  has a strategy  $\sigma_S(\theta) : \Theta \rightarrow \prod_{i \in \mathcal{N}} \mathcal{R}^i$ . Receiver  $R^i$  has a strategy  $\sigma_R^i(r^i) : \mathcal{R}^i \rightarrow \mathcal{A}$ . The strategies of both  $S$  and  $R^i$  induce distributions  $\rho_S(\mathbf{r}|\theta) = \{\rho_S^1(r^1|\theta), \dots, \rho_S^n(r^n|\theta)\}$  and  $\rho_R^i(a^i|r^i)$ , respectively. A pair of distributions  $(\rho_S^i(r^i|\theta), \rho_R^i(a^i|r^i))$  satisfies

$$\int_{r^i \in \mathcal{R}^i} \rho_S^i(r^i|\theta) dr^i = 1, \forall i \in \mathcal{N}, \forall \theta \in \Theta,$$

$$\sum_{a^i \in \mathcal{A}} \rho_R^i(a^i|r^i) = 1, \forall i \in \mathcal{N}, \forall r^i \in \mathcal{R}^i.$$

In this game,  $R^i$ , who observes the message  $r^i$  before choosing  $a^i$ , updates its beliefs about the type  $\theta$  using the Bayes' rule. The goal of the  $R^i$  is to choose an action  $a^i$  to minimize its expected cost  $c_R^i$  given a posterior belief  $\mu^i(\theta|r^i)$ , while the goal of the sender is to choose a signal  $r^i$  to minimize the cost  $c_S$  by anticipating the behavior of the receiver  $R^i$ . The following definition identifies an PBNE, where both the  $S$  and  $R^i$  choose their best strategy to minimize the their cost function.

**Definition 1:** A PBNE of the signaling game is a strategy profile  $\{\sigma_S, \sigma_R^i, i \in \mathcal{N}\}$  and posterior beliefs  $\mu^i(\theta|r^i)$ , where  $i \in \mathcal{N}$ , such that:

$$\forall \theta, \sigma_R^i(r^i) \in \arg \min_{a^i \in \mathcal{A}} \sum_{\theta} \mu^i(\theta|r^i) c_R^i(r^i, a^i, \theta), \quad (2)$$

$$\forall \theta, \sigma_S(\theta) \in \arg \min_{\mathbf{r} \in \mathcal{R}} c_S(\mathbf{r}, \sigma_R, \theta), \quad (3)$$

and

$$\mu^i(\theta|r^i) = \begin{cases} \frac{p^i(\theta) \rho_S^i(r^i|\theta)}{\sum_{\theta' \in \Theta} \pi_S(\theta', r^i)}, & \text{if } \sum_{\theta' \in \Theta} \pi_S(\theta', r^i) > 0, \\ \text{any distribution,} & \text{if } \sum_{\theta' \in \Theta} \pi_S(\theta', r^i) = 0. \end{cases}$$

where  $\mathcal{R} = \prod_{i \in \mathcal{N}} \mathcal{R}^i$ ,  $\sigma_S(\theta) = (\sigma_S^1(\theta), \dots, \sigma_S^n(\theta))$ , and  $\pi_S(\theta', r^i) = p^i(\theta') \rho_S^i(r^i|\theta')$ .

**Remark 1:** There are two important classes of equilibria in a signaling game. The first one is the pooling equilibrium, stating that different types of senders choose the same strategy. The second one is the separating equilibrium, stating that different types of senders choose distinct strategies.

**Remark 2:** The sender's problem can be divided into  $n$  sub-problems if the interactions between  $S$  and  $R^i$  is independent of the strategies of other ASs. The decoupling process is given by

$$\min_{\mathbf{r} \in \mathcal{R}} c_S(\mathbf{r}, \sigma_R, \theta)$$

$$= \min_{\mathbf{r} \in \mathcal{R}} \sum_{i \in \mathcal{N}} c_S^i(r^i, \sigma_R^i, \theta) = \sum_{i \in \mathcal{N}} \min_{r^i \in \mathcal{R}^i} c_S^i(r^i, \sigma_R^i, \theta),$$

where  $c_S^i(r^i, \sigma_R^i, \theta)$  is the sub-problem for  $S$ . Therefore, instead of solving (3),  $S$  can generate  $r_k^i$  by solving  $n$  sub-problems,

$$\sigma_S^i \in \arg \min_{r^i \in \mathcal{R}^i} c_S^i(r^i, \sigma_R^i, \theta).$$

### D. Detection Range and Reference Set

Each AS has a range of detection of neighboring systems. Let the distance  $d(z_k^i, z_k^j) \triangleq \|z_k^i - z_k^j\|$ . To avoid collision, we assume that each  $R^i$  is equipped with a special local sensor which can sense objects within a limited range  $L_d$ , i.e., AS  $i$  can sense AS  $j$ 's position at time  $k$  if  $\|z_k^i - z_k^j\| \leq L_d^i$ . Given the detection range  $L_d^i$ , we define two sets as follows:

$$\mathcal{D}_k^i \triangleq \{z_k^j : d(z_k^i, z_k^j) \leq L_d^i\},$$

where  $\mathcal{D}_k^i$  is the detection set of  $R^i$  at time  $k$ .

To avoid collision, we define the reference set by  $\mathcal{R}_k^i = \{r_k^i : d(r_k^i, z_k^i) \leq L_r^i\}$ , where  $L_r^i$  is a reference radius, which is a constant scalar. The radius  $L_r^i$  is design to be  $L_r^i \leq L_d^i$ . Therefore,  $R^i$  can sense all the objects inside  $\mathcal{R}_k^i$  at time  $k$ .

### III. ANALYSIS OF THE CYBER-PHYSICAL SIGNALING GAME

In this section, we will define the cost functions of  $S$  and  $R^i$ , and analyze the PBNE strategies of the game.

#### A. Strategy and Cost Function

The strategies of  $S$  and  $R^i$  depend on the type  $\theta$  and message  $r^i$ , respectively. We assume that the normal sender ( $\theta = \theta_0$ ) sends a legitimate reference signal  $\bar{r}_k^i$  for  $R^i$ . Hence a penalty will incur when a reference signal  $r_k^i$  deviates from  $\bar{r}_k^i$ . When  $\theta = \theta_1$ , a malicious sender attempts to send a reference signal  $r_k^i$  to collide  $R^i$  with a static obstacle or other ASs. Without loss of generality, we assume that the obstacle is AS  $j$  whose position is given by  $z_k^j$ . For a static object,  $z_k^j$  will be independent of  $j$  and  $k$ . To achieve collision, the attacker aims to minimize the distance between  $r_k^i$  and  $z_k^j$ .

To simply the problem, we assume that the interaction between  $S$  and  $R^i$  is independent of other receivers, hence the problem of  $S$  can be divided into  $n$  sub-problems as pointed out in Remark 2. Next, we define the cost functions for  $S$  and  $R^i$ .

1) *The cost function of  $R^i$* : The cost function  $c_R^i$  for  $R^i$  depends on the type of  $S$ , the actions of the sender  $r^i$  and the action  $a^i$  of the receiver. When  $S$  is normal ( $\theta = \theta_0$ ), we let  $J_{\theta_0} : \mathcal{A}_i \times \mathcal{R}^i \rightarrow \mathcal{R}$  be the cost function for  $R^i$ :

$$J_{\theta_0}(a_k^i, r_k^i) \triangleq (1 - a_k^i) \left( J_c(x_k^i, \bar{u}_k^i, z_k^i) + \phi^i \right), \quad (4)$$

where  $\bar{u}_k^i = (\bar{u}_{k|k}^i, \dots, \bar{u}_{k+N-1|k}^i)$  is the control input for  $R^i$  to stop at point  $z_k^i$ ,  $\phi^i \in \mathbb{R}_+$  is a mission cost when  $R^i$  does not follow the reference  $r_k^i$ . The cost function (4) shows that if the sender is normal and the  $R^i$  chooses  $a_k^i = 1$  to fulfill the mission,  $R^i$  will have a zero cost. Otherwise,  $R^i$  has to pay a cost for staying at position  $z_k^i$  and a mission cost  $\phi^i$ . To simplify the problem, we assume that the mission cost  $\phi^i$  is a constant.

When  $S$  is malicious ( $\theta = \theta_1$ ), we denote by  $J_{\theta_1} : \mathcal{A}_i \times \mathcal{R}^i \rightarrow \mathbb{R}$  the cost for the distance between the reference  $r_k^i$

and the position  $z_k^j$  of  $R^j$  when  $z_k^j \in \mathcal{D}_k^i$ . Hence,  $J_{\theta_1}$  is given by

$$J_{\theta_1}(a_k^i, r_k^i) \triangleq \frac{a_k^i K^i}{d(r_k^i, z_k^j)} + \frac{(1 - a_k^i) K^i}{d(z_k^i, z_k^j)}, \quad (5)$$

where  $K^i \in \mathbb{R}_+$  is a tuning parameter. Give the cost (4) and (5), the expected cost function of  $R^i$  is as follows:

$$\mathbb{E}[c_R^i(a_k^i, r_k^i, \theta)] \triangleq \sum_{\theta \in \Theta} \mu^i(\theta | r^i) J_{\theta}(a_k^i, r_k^i), \quad (6)$$

2) *The cost function of  $S$  with its type  $\theta_0$* : The sender  $S$  has two types. When  $S$  is normal, the cost function  $c_S^i : \mathcal{A}_i \times \mathcal{R}^i \times \Theta \rightarrow \mathbb{R}$  of  $S$  is defined by

$$c_S^i(a_k^i, r_k^i, \theta_0) \triangleq a_k^i d(r_k^i, \bar{r}_k^i) + (1 - a_k^i) d(z_k^i, \bar{r}_k^i). \quad (7)$$

The cost function (7) shows that when  $R^i$  chooses  $a^i = 1$ , the normal sender aims to minimize  $d(r_k^i, \bar{r}_k^i)$ . Otherwise, it will achieve  $d(z_k^i, \bar{r}_k^i)$  when  $a_k^i = 0$ .

3) *The cost function of  $S$  with its type  $\theta_1$* : When  $S$  is normal, the cost function  $c_S^i : \mathcal{A}_i \times \mathcal{R}^i \times \Theta \rightarrow \mathcal{R}$  of  $S$  is defined by

$$c_S^i(a_k^i, r_k^i, \theta_1) \triangleq a_k^i d(r_k^i, z_k^j) + (1 - a_k^i) d(z_k^i, z_k^j). \quad (8)$$

This cost function illustrates that when  $a_k^i = 1$ , the malicious sender aims to minimize  $d(r_k^i, z_k^j)$  to cause a collision between  $R^i$  and  $R^j$  or a static object if  $z_k^j$  is taken as a constant.

#### B. Signaling Game Analysis

To characterize the equilibrium for this cyber-physical game, we first investigate the strategy of receiver  $R^i$ .

*Lemma 1*: Let  $R^i$ 's cost function defined by (6). A pure strategy of  $R^i$  at time  $k$  is given by

$$a_k^i = \mathbf{1}_{\{r_k^i \notin \Omega_k^{ij}(\mu_k^i(\theta | r_k^i))\}}, \quad (9)$$

where

$$\Omega_k^{ij}(\mu_k^i(\theta | r_k^i)) = \{r_k^i : d(r_k^i, z_k^j) < \eta_k^{ij}(\mu_k^i(\theta | r_k^i))\},$$

$$\eta_k^{ij}(\mu_k^i(\theta | r_k^i)) \triangleq$$

$$\frac{\mu_k^i(\theta_1 | r_k^i) K^i d(z_k^i, z_k^j)}{\mu_k^i(\theta_1 | r_k^i) K^i + \mu_k^i(\theta_0 | r_k^i) J_{\theta_0}(a_k^i = 1, r_k^i) d(z_k^i, z_k^j)}.$$

*Remark 3*: The  $R^i$ 's strategy given by Lemma 1 shows that if  $r_k^i \in \Omega_k^{ij}(\mu_k^i(\theta | r_k^i))$ , then  $R^i$  will reject  $r_k^i$ . This strategy can protect  $R^i$  from collision with  $R^j$  as it keeps  $R^i$  from  $R^j$  with a safe distance  $\eta_k^{ij}(\mu_k^i(\theta | r_k^i))$ . Here, we call  $\Omega_k^{ij}$  a danger zone.

*Corollary 1*: Let  $\mathcal{M}_k^i = \{j : z_k^j \in \mathcal{D}_k^i\}$  be an index set of AS within the detection range  $L_d^i$ . If the size  $|\mathcal{M}_k^i| \geq 2$ , then, the strategy of  $R^i$  is given by

$$a_k^i = \mathbf{1}_{\{r_k^i \notin \bigcup_{j \in \mathcal{M}_k^i} \Omega_k^{ij}(\mu_k^i(\theta | r_k^i))\}}. \quad (10)$$

*Remark 4*: Corollary 1 shows that if more than one AS belong to the set  $\mathcal{D}_k^i$ , then, the danger zone of  $R^i$  is the joint of the sub-danger zone generated by  $R^j$ , where  $j \in \mathcal{M}_k^i$ . Fig. 2 shows an example that three ASs are inside  $\mathcal{D}_k^i$  of  $R^i$ .  $R^i$  rejects the  $r_k^i$  if it is inside the gray circles given in Fig. 2.

**Lemma 3:** When  $\theta = \theta_0$ , the strategy of  $S$  depends on  $\bar{r}_k^i$ . When  $d(\bar{r}_k^i, z_k^j) > \eta_k^{ij}(\mu^i(\theta|r_k^i))$ , the strategy of  $S$  is given by  $\sigma_S^i(\theta_0) = \bar{r}_k^i$ . When  $d(\bar{r}_k^i, z_k^j) \leq \eta_k^{ij}$ , the strategy of  $S$  is given by  $\sigma_S^i(\theta_0) = r^i(\theta_0)$ , where  $r^i(\theta_0)$  corresponds to the red dot in Fig. 3(b).

**Lemma 3:** When  $\theta = \theta_0$ , the strategy of  $S$  depends on  $\bar{r}_k^i$ . When  $d(\bar{r}_k^i, z_k^j) > \eta_k^{ij}(\mu^i(\theta|r_k^i))$ , the strategy of  $S$  is given by  $\sigma_S^i(\theta_0) = \bar{r}_k^i$ . When  $d(\bar{r}_k^i, z_k^j) \leq \eta_k^{ij}$ , the strategy of  $S$  is given by  $\sigma_S^i(\theta_0) = r^i(\theta_0)$ , where  $r^i(\theta_0)$  corresponds to the red dot in Fig. 3(b).

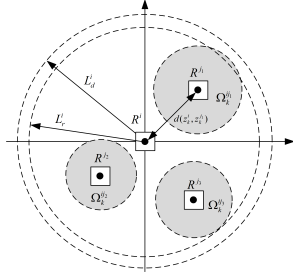


Fig. 2. The strategy of the receiver: Having the knowledge of the beliefs  $\mu_k^i(\theta|r^i)$  and the position  $z_k^j$ ,  $R^i$  computes danger zones  $\Omega_k^{ij}$ ;  $R^i$  only accepts  $r_k^i$  if  $r_k^i \notin \bigcup_j \Omega_k^{ij}(\mu_k^i, k)$ .

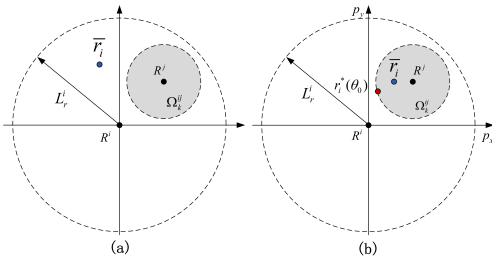


Fig. 3. The best response of  $S$  of type  $\theta_0$ : If  $d(\bar{r}_k^i, z_k^j) > \eta_k^{ij}(\mu^i(\theta|r_k^i))$ , then the best response of  $S$  is  $\bar{r}_k^i$ ; if  $d(\bar{r}_k^i, z_k^j) \leq \eta_k^{ij}(\mu^i(\theta|r_k^i))$ , then the best response of  $S$  is  $r^i(\theta_0)$  (the red dot).

With Lemmas 1, 2 and 3, the following Theorem 1 characterizes the PBNEs in our signal game after observing message  $r_k^i$  of the sender.

**Theorem 1:** Two PBNEs exist in this cyber-physical signaling game: One is a separating equilibrium, and the other is a pooling equilibrium. They are presented as follows:

(1)  $\forall j \in \mathcal{M}_k^i$ , if  $d(\bar{r}_k^i, z_k^j) > d(z_k^i, z_k^j)$ , a separating equilibrium exists, which is given by

$$\sigma_R^i(r_k^i) = \mathbf{1}_{\{r_k^i \notin \Omega_k^i\}}, \quad (11)$$

$$\sigma_S^i(\theta, k) = \begin{cases} z_k^i, & \text{if } \theta = \theta_1 \\ \bar{r}_k^i, & \text{if } \theta = \theta_0 \end{cases}, \quad (12)$$

$$\mu_k^i(\theta_1|r_k^i) = 1, \quad \forall r^i \in \Omega_k^i, \quad (13)$$

$$\mu_k^i(\theta_0|r_k^i) = 1, \quad \forall r^i \in \bar{\Omega}_k^i \cap \mathcal{R}_k^i, \quad (14)$$

where  $\bar{\Omega}_k^i$  is the complementary set of  $\Omega_k^i$ , and

$$\begin{aligned}\mathbf{\Omega}_k^i &= \bigcup_{j \in \mathcal{M}_k^i} \hat{\Omega}_k^{ij}, \\ \hat{\Omega}_k^{ij} &= \{r_k^i : d(r_k^i, z_k^j) < d(z_k^i, z_k^j)\}.\end{aligned}\tag{15}$$

$$\sigma_R^i(r_k^i) = \mathbf{1}_{\{r_k^i \notin \Omega_k^i\}}, \quad (16)$$

$$\sigma_S^i(\theta, k) = z_k^i, \quad \forall \theta \in \Theta, \quad (17)$$

$$\mu_k^i(\theta_1 | r_k^i) = 1, \quad \forall r^i \in \Omega_k^i, \quad (18)$$

$$\mu_k^i(\theta_0|r_k^i) = 1, \quad \forall r^i \in \bar{\Omega}_k^i \cap \mathcal{R}_k^i, \quad (19)$$

where  $\Omega_k^i$  is defined by (15).

*Remark 5:* Both the equilibria presented in Theorem 1 can protect each system from collision as the equilibria can guarantee that  $R^i$  only accepts  $r_k^i$  when  $d(r_k^i, z_k^j) \geq d(r_k^i, z_k^j)$  for  $j \in \mathcal{M}_k^i$ . However,  $R^i$  cannot distinguish the type of  $S$  when the strategies of  $S$  and  $R^i$  arrive at the pooling equilibrium.

The following corollary presents a sufficient condition for the existence of a unique separating PBNE.

*Corollary 2:* If  $d(\bar{r}_k^i, z_k^j) > d(z_k^i, z_k^j)$  for all  $j \in \mathcal{M}_k^i$ , a unique separating PBNE defined by (11) and (12) with the beliefs (13), and (14) exists in the signaling game.

*Remark 6:* Corollary 2 shows that if the normal sender chooses a safe  $\bar{r}^i$  that satisfies  $d(\bar{r}_k^i, z_k^j) > d(z_k^i, z_k^j), \forall j \in \mathcal{M}_k^i$ , then a unique separating PBNE exists. This PBNE benefits the whole system as  $R^i$  only accepts the reference from the normal sender.

## IV. SIMULATION RESULTS

A UAV example is used in the experiments. The program of the UAV and the control station are performed on the different blocks in the Simulink of MATLAB 2014b, running on the workstation with an Intel Core i7-4770 processor and a 12-G RAM. All UAVs are assumed to fly at the same altitude.

Before testing the proposed mechanism, we first show the impact of the MITM attacks. Fig. 4 illustrates the scenario of Suicidal Attack, where the attacker gives a fake reference to lead all the UAVs to crash into a building (the gray block in Fig. 4). In the second experiment, the attacker indicates a point to assemble all the UAVs to launch a collision attack. Consequently, all the UAVs collide at the point (60, 60) in Fig. 5.

The next part of the simulation is to validate the mechanism under the two MITM attack models given in Section 2.2. Fig. 6 shows that the attacker aims to deviate all the UAVs to crash to the building. However, UAV<sub>1</sub> chooses to stop moving when it approaches the building. The other UAVs following UAV<sub>1</sub> stop and form a line with two ASs separated by a safe distance. The second plot in Fig. 7 shows the posterior belief  $\mu^1(\theta_1|r^1)$  of UAV<sub>1</sub> converges to 1, indicating that UAV<sub>1</sub> learns the occurrence of the attack, and it chooses  $a^1 = 0$  at  $k = 170$  to avoid colliding with the building. Fig. 8 shows that the distance between the UAV<sub>1</sub> and the building under the suicidal attack. With proposed mechanism, the distance does not decay to 0. These results show that all the UAVs avoid collisions under the suicidal attack. Finally, we test our mechanism under the collision

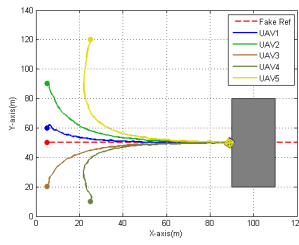


Fig. 4. An example of suicidal attack: The attacker generates a fake reference (red dash line) to lead all the UAVs to crash to a building.

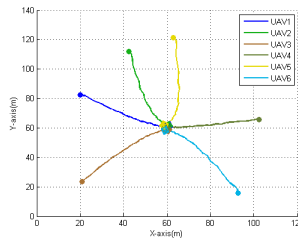


Fig. 5. An example of collision attack for UAVs: Attacker sends a reference point to collide all the UAVs at one point.

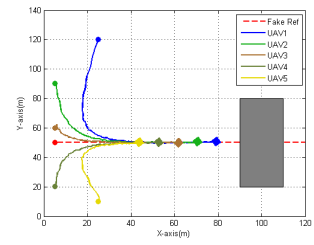


Fig. 6. The UAVs' system with defense strategy: All the UAVs stop in a line to avoid collisions.

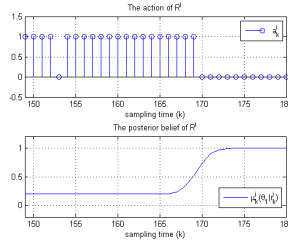


Fig. 7. Action  $\alpha^1$  and the posterior beliefs  $\mu^1(\theta_1|r^1)$  for UAV<sub>1</sub>: When an attack occurs at  $k = 166$ , UAV<sub>1</sub> updates  $\mu^1(\theta_1|r^1)$ , which converges to 1 when UAV<sub>1</sub> approaches the building. At  $k = 170$ , UAV<sub>1</sub> chooses  $\alpha^1 = 0$ .

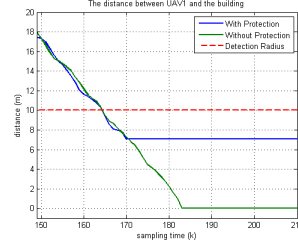


Fig. 8. The distance between the UAV<sub>1</sub> and the building under the suicidal attack: (1) without protection, the distance decays to 0; (2) with proposed mechanism, the distance decays to a positive constant.

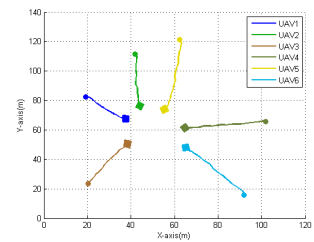


Fig. 9. With the proposed mechanism, under the collision attack, each UAV avoids collisions by choosing hovering with a safe distance from other UAVs.

attack. Fig. 9 shows that all the UAVs avoid collisions, and stop within a safe distance from the target.

## V. CONCLUSIONS

Autonomous systems (ASs) are increasingly vulnerable to cyber attacks. The cyber-physical design of security mechanism for ASs is critical to provide performance guarantee and minimize the damage caused by sophisticated attacks. In this paper, we have proposed an integrated game-theoretic framework to address the cyber-physical nature and real-time requirements of the ASs. We have used signaling games and minimax model predictive control methods to develop co-design methodologies for impact-aware cyber defense mechanisms and the cyber-aware defense strategies. We have shown that under a sufficient condition, a unique PBNE exists in our game model, protecting the ASs from collisions.

## REFERENCES

- [1] "Amazon Prime Air," <http://www.amazon.com/b?node=8037720011> [Last accessed in March 4, 2015].
- [2] E. Guizzo, "How googles self-driving car works," *IEEE Spectrum Online*, October, vol. 18, 2011.
- [3] Y. Tipsuwan and M.-Y. Chow, "Control methodologies in networked control systems," *Control engineering practice*, vol. 11, no. 10, pp. 1099–1111, 2003.
- [4] J. P. Hespanha, P. Naghshtabrizi, and Y. Xu, "A survey of recent results in networked control systems," *Proceedings of the IEEE*, vol. 95, no. 1, p. 138, 2007.
- [5] R. Langner, "Stuxnet: Dissecting a cyberwarfare weapon," *Security & Privacy, IEEE*, vol. 9, no. 3, pp. 49–51, 2011.
- [6] B. Bencsáth, G. Pék, L. Buttyán, and M. Félégyházi, "Duqu: Analysis, detection, and lessons learned," in *ACM European Workshop on System Security (EuroSec)*, vol. 2012, 2012.
- [7] C. Tankard, "Advanced persistent threats and how to monitor and deter them," *Network security*, vol. 2011, no. 8, pp. 16–19, 2011.
- [8] A. J. Kerns, D. P. Shepard, J. A. Bhatti, and T. E. Humphreys, "Unmanned aircraft capture and control via gps spoofing," *Journal of Field Robotics*, vol. 31, no. 4, pp. 617–636, 2014.
- [9] Y. Desmedt, "Man-in-the-middle attack," in *Encyclopedia of Cryptography and Security*. Springer, 2011, pp. 759–759.
- [10] J.-S. Kim, T.-W. Yoon, A. Jadbabaie, and C. De Persis, "Input-to-state stabilizing MPC for neutrally stable linear systems subject to input constraints," in *43rd IEEE Conference on Decision and Control (CDC)*, vol. 5, 2004, pp. 5041–5046.
- [11] C. E. Garcia, D. M. Prett, and M. Morari, "Model predictive control: theory and practice a survey," *Automatica*, vol. 25, no. 3, pp. 335–348, 1989.
- [12] A. Teixeira, D. Pérez, H. Sandberg, and K. H. Johansson, "Attack models and scenarios for networked control systems," in *Proceedings of the 1st international conference on High Confidence Networked Systems*. ACM, 2012, pp. 55–64.
- [13] S. Bhattacharya and T. Basar, "Game-theoretic analysis of an aerial jamming attack on a UAV communication network," in *American Control Conference (ACC), 2010*. IEEE, 2010, pp. 818–823.
- [14] W. Gao, T. Morris, B. Reaves, and D. Richey, "On SCADA control system command and response injection and intrusion detection," in *eCrime Researchers Summit (eCrime), 2010*. IEEE, 2010, pp. 1–9.
- [15] Q. Zhu, L. Bushnell, and T. Başar, "Game-theoretic analysis of node capture and cloning attack with multiple attackers in wireless sensor networks," in *Decision and Control (CDC), 2012 IEEE 51st Annual Conference on*. IEEE, 2012, pp. 3404–3411.
- [16] D. Welch and S. Lathrop, "Wireless security threat taxonomy," in *Information Assurance Workshop, 2003. IEEE Systems, Man and Cybernetics Society*. IEEE, 2003, pp. 76–83.
- [17] J. Zhen and S. Srinivas, "Preventing replay attacks for secure routing in ad hoc networks," in *Ad-Hoc, Mobile, and Wireless Networks*. Springer, 2003, pp. 140–150.
- [18] B. Kannhavong, H. Nakayama, Y. Nemoto, N. Kato, and A. Jamalipour, "A survey of routing attacks in mobile ad hoc networks," *Wireless communications, IEEE*, vol. 14, no. 5, pp. 85–91, 2007.
- [19] L. Xie and S. Zhu, "Message dropping attacks in overlay networks: Attack detection and attacker identification," *ACM Transactions on Information and System Security (TISSEC)*, vol. 11, no. 3, p. 15, 2008.
- [20] M. H. Manshaei, Q. Zhu, T. Alpcan, T. Başar, and J.-P. Hubaux, "Game theory meets network security and privacy," *ACM Computing Surveys (CSUR)*, vol. 45, no. 3, p. 25, 2013.