

# A zonotopic set-invariance analysis of replay attacks affecting the supervisory layer

Carlos Trapiello<sup>a,b,\*</sup>, Vicenç Puig<sup>a,b</sup>, Damiano Rotondo<sup>c</sup>

<sup>a</sup>Advanced Control Systems Group, Automatic Control Department, Universitat Politècnica de Catalunya (UPC), Rambla Sant Nebridi 10, Terrassa 08222, Spain

<sup>b</sup>Institut de Robòtica i Informàtica Industrial, CSIC-UPC, Llorens i Artigas 4-6, Barcelona 08028, Spain

<sup>c</sup>Department of Electrical and Computer Engineering (IDE), University of Stavanger (UiS), Kristine Bonnevis vei 22, 4021 Stavanger, Norway

---

## Abstract

This paper presents a zonotopic set-invariance analysis of replay attacks affecting the communication network that serves the supervisory layer of complex control systems using an observer-based detection scheme. Depending on the attacker's access to the system's resources, two scenarios are considered: I) Sensors and controller data are counterfeited; II) Only sensor measurements are counterfeited. The effect of a physical attack against the plant during the data replay is also taken into consideration. The representation of invariant sets as zonotopes allows to derive analytical expressions for attack detectability under the presence of bounded uncertainties. The validity of the analysis is demonstrated through simulations using a quadruple-tank process.

**Keywords:** Invariant sets, replay attack, zonotopes, attack detection.

---

## 1. Introduction

The migration from traditional point-to-point control schemes to widely interconnected systems, in conjunction with an increasing number of registered attacks [1], has awakened the interest on the study of cyber attacks on control systems. Consequently, the so-called *secure control* has found increasing interest since the end of the last century. In this regard, works like [2, 3] propose a general framework for the analysis of security on networked control systems, while other works like [4] have modelled the effect of cyber attacks compromising measurement and actuator data integrity on the physical dynamics of a system.

One of the main threats to control systems are stealthy attacks, i.e., attacks in which the attacker aims at remaining undetected by anomalies detectors [5]. In order to achieve undetectability, the attacker must be able to feed data that is consistent with the nominal system operation to the monitoring system. Among the different stealthy attacks reported in the literature, there are: false data injection attacks [6], zero dynamics attacks [7], covert attacks [8] and replay attacks [9]. This paper focuses on the latter type of attacks.

The vast majority of works related to replay attacks takes into consideration that the control loop is closed remotely by means of a communications network which is prone to cyber attacks. Accordingly, the standard replay attack formulation considers that the control loop has been affected either by deceiving the system controller

with previously recorded measurements [9, 10, 11, 12], or by replaying back previous control actions directly to the system [13, 14]. Nevertheless, in many industrial control systems the low-level controller (regulatory layer) makes use of dedicated networks which are hard to access while, on the other hand, it is common that the supervisory layer, in charge of the system monitoring and set-point reference generation, operates remotely. This difference has been taken into account, for example, in the study of cyber-attacks affecting the load frequency control of power systems [15]. In this line, this paper is devised from the supervisor's point of view, analysing different replay attack scenarios that arise when the system operation is assessed using a state estimator located at the supervisory layer.

Moreover, a high percentage of security-related works (not only concerning attacks but also faults) makes use of assumptions about the statistical properties of the uncertainties. A different approach is the use of set-theoretic methods, which are built upon norm-bounded uncertainty assumptions. These techniques have proven useful in fault-related secure control, as they allow to construct sets for the system in healthy and faulty operations, so that it becomes possible to infer deterministically whether a system is working under nominal or faulty conditions [16]. In this regard, ellipsoidal sets have been employed in [17] to propose different security metrics, and design secure controllers, against stealthy attacks affecting the control loop. On the other hand, set-based attack detectors have been used in [18] to detect bias injection attacks in a networked power plant, and in [19, 20] combined altogether with set-theoretic controllers in an attack resilient control scheme. However, these latter detectors do not rely on an

---

\*Corresponding author

Email address: carlos.trapiello@upc.edu (Carlos Trapiello)

observer, and thus are subject to some limiting assumptions concerning the knowledge of the initial state of the system.

Within the set-theoretic techniques, positive set invariance [21] is a common analysis tool for systems affected by bounded disturbances. By computing the healthy/attacked residual invariant sets, it can be established that whenever the residual vector exits the healthy set, attack detection is achieved. Note that the inherent steady-state conditions of replay attacks suit perfectly with set invariant tools. Among the possible set representations, zonotopes will be employed in this paper due to their flexibility and the capability to separate the center (nominal) evolution from the uncertainty evolution in the generators representation [22].

The motivation of the paper is characterizing the detectability of the replay attacks in a guaranteed manner using set-invariance analysis when an observer-based detection scheme is used. Therefore, system's vulnerabilities are analysed for the case in which a malicious attacker compromises the communication between the regulatory and supervisory layers of a control system, while the low-level controller at the regulatory layer remains unaffected. The following scenarios are considered: I) the attacker is able to record and replay sensors and controller data, causing the supervisory layer to operate based on false data; II) the attacker is able to replay sensors data to the supervisory layer while the low-level control actions are received unaffected. This scenario models the case where the supervisory layer operates based on a set of sensors installed for system monitoring and that may differ from the sensors used for control. For both scenarios, a positive invariance approach is developed for analysing attack detectability with respect to the set-point reference signal injected from the supervisory layer and an external attack conducted over the plant, where the invariant sets are represented as zonotopes. Conditions to guarantee attack detectability are derived under the assumption that the disturbances are bounded.

The article is structured as follows: Section 2 introduces some preliminary developments regarding the zonotopic representation of invariant sets, as well as the different phases that constitute the attack. Section 3 is devoted to the description of the system in healthy operation, while the assumptions on the system operation during the record phase are detailed in Section 4. Section 5 presents the analysis of the detectability of the attack during the replay phase. An illustrative example is presented in Section 6. Finally, Section 7 presents concluding remarks.

## 2. Preliminaries

### 2.1. Zonotopes and basic set operations

Zonotopes are centrally symmetric convex polytopes that can be described as Minkowski sums of line segments [23]. In the generator representation, a zonotope

$\mathcal{Z}$  is described by its center  $c \in \mathbb{R}^n$  and generators  $g_1, \dots, g_m \in \mathbb{R}^n$  as  $\mathcal{Z} = \{c + H\xi : \xi \in \mathbb{R}^m, \|\xi\|_\infty \leq 1\}$  where  $H \equiv [g_1, \dots, g_m]$  indicates the generators matrix and the ratio  $m/n$  is the order of the zonotope. For simplicity, zonotopes will be denoted by  $\mathcal{Z} = \langle c, H \rangle$ .

Let the sets  $\mathcal{Z}, \mathcal{W} \subset \mathbb{R}^n$ , the matrix  $P \in \mathbb{R}^{k \times n}$ , and define

$$P\mathcal{Z} \equiv \{Pz : z \in \mathcal{Z}\}, \quad (1a)$$

$$\mathcal{Z} \oplus \mathcal{W} \equiv \{z + w : z \in \mathcal{Z}, w \in \mathcal{W}\}. \quad (1b)$$

Zonotopes are closed under previous set operations, i.e., when  $\mathcal{Z} = \langle c_z, H_z \rangle$  and  $\mathcal{W} = \langle c_w, H_w \rangle$  are zonotopes, linear mappings (1a) and Minkowski sums (1b) are also zonotopes which can be computed as

$$P\mathcal{Z} = \langle Pc_z, PH_z \rangle, \quad (2a)$$

$$\mathcal{Z} \oplus \mathcal{W} = \langle c_z + c_w, [H_z \ H_w] \rangle. \quad (2b)$$

### 2.2. Invariant sets

Let us define a discrete-time linear time-invariant (LTI) system

$$x^+ = Ax + B\delta, \quad (3)$$

where  $x \in \mathbb{R}^n$  is the system state,  $x^+ \in \mathbb{R}^n$  its successor and  $\delta \in \mathbb{R}^{n_\delta}$  is a disturbance constrained to the compact zonotopic set  $\Delta = \langle c_\Delta, H_\Delta \rangle \subset \mathbb{R}^{n_\delta}$ . Besides,  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times n_\delta}$  are constant matrices with  $A$  an asymptotically stable matrix (all the eigenvalues of  $A$  are strictly inside the unit disk).

**Definition 1** (Robust positive invariance). The set  $\Omega \subset \mathbb{R}^n$  is said to be *robustly positively invariant* (RPI) for the system (3) and disturbance set  $\Delta$ , if  $Ax + B\delta \in \Omega$  for all  $x \in \Omega$  and all  $\delta \in \Delta$ . Equivalently,  $\Omega$  is RPI if and only if  $A\Omega \oplus B\Delta \subseteq \Omega$ .

**Definition 2** (Minimal RPI). The minimal RPI (mRPI) set of (3) is the RPI set in  $\mathbb{R}^n$  that is contained in every closed RPI set of (3) and disturbance set  $\Delta$ .

For LTI asymptotically stable systems like (3), the mRPI set exists and is unique and compact [24, Sec. IV]. In addition, such mRPI set is the limit set of all trajectories of the system. Henceforth, a zonotopic  $\epsilon$ -approximation of the mRPI set will be computed by means of Algorithm 1 detailed in Appendix A. When referring to Algorithm 1, the center and the generators matrix recursion will be provided.

For a detailed analysis on set invariance the reader is referred to comprehensive studies [24, 25].

### 2.3. Attack time windows

In the attack under study, it is assumed that at a first stage a malicious attacker secretly records the data transmitted from the regulatory layer to the supervisory layer. Then, the recorded data are replayed back to the supervisory layer with the intention of masking a physical attack conducted over the plant. Consequently, the following time windows are defined:

1. **Record window:** transmitted data are assumed to be recorded for  $\mathcal{K}_{rec} = \{k \in \mathbb{N} : k \in [k_0, k_0 + l - 1]\}$ , where  $l \in \mathbb{N}$  denotes the size of the record window.
2. **Replay window:** real data are replaced for  $\mathcal{K}_{rep} = \{k \in \mathbb{N} : k \in [k_1 + (n - 1)l, k_1 + nl - 1], \forall n \in \{1, \dots, n_r\}\}$ , where  $n_r \in \mathbb{N}^+$  accounts for the total number of repetitions of the recorded sequence.
3. **Physical attack window:** a physical attack against the plant is launched for  $\mathcal{K}_{phy} = \{k \in \mathbb{N} : k \in [k_2, k_3]\} \subseteq \mathcal{K}_{rep}$ , i.e.,  $k_2 \geq k_1$  and  $k_3 \leq k_1 + n_r l - 1$ .

Whenever one of the above temporal sets is mentioned, it is assumed implicitly that the index  $k$  lies within it.

### 3. System under healthy operation

Let us consider the following system

$$\begin{aligned} x_{k+1} &= Ax_k + Bu_k + E_w w_k, \\ y_k &= Cx_k + E_v v_k, \end{aligned} \quad (4)$$

where  $x_k \in \mathbb{R}^{n_x}$  is the state vector,  $u_k \in \mathbb{R}^{n_u}$  is the applied control action and  $y_k \in \mathbb{R}^{n_y}$  corresponds to the sensor measurements at time instant  $k$ . Furthermore,  $w_k \in \mathbb{R}^{n_w}$  and  $v_k \in \mathbb{R}^{n_v}$  represent the process disturbances and measurement noise, respectively. Henceforth, the index  $k + 1$  will be replaced by the superscript  $+$  and  $k$  will be omitted for the sake of simplified notations.

**Assumption 1.** The pair  $(A, B)$  is asymptotically stabilizable and the pair  $(A, C)$  is asymptotically detectable.

**Assumption 2.** Uncertainties are bounded by

$$w \in \mathcal{W} = \langle c_w, H_w \rangle, \quad v \in \mathcal{V} = \langle c_v, H_v \rangle, \quad (5)$$

with  $H_w \in \mathbb{R}^{n_w \times m_w}$ ,  $H_v \in \mathbb{R}^{n_v \times m_v}$  known generator matrices and  $c_w \in \mathbb{R}^{n_w}$ ,  $c_v \in \mathbb{R}^{n_v}$  known zonotope centers.

The control objective is to regulate the plant tracking error defined at each sample as  $z = x - x_{ref}$ , where  $x_{ref} \in \mathbb{R}^{n_x}$  is the reference signal governed by

$$x_{ref}^+ = Ax_{ref} + Bu_{ref}, \quad (6)$$

and  $u_{ref} \in \mathbb{R}^{n_u}$  is the reference signal generated at the supervisory layer, which is sent to the low-level controller (regulatory layer) in charge of regulating the plant's tracking error (see Figure 1). For a desired output set-point  $y_{ref} = Cx_{ref}$ , the corresponding  $u_{ref}$  signal can be obtained by means of classical model inversion-based feed-forward schemes [26].

#### 3.1. Regulatory layer

In order to satisfy the control objective, the low-level controller is assumed to perform an estimate-feedback control action based on the estimates provided by its own set-based observer [27, 28, 29], and that differs from the one used in the supervisory layer. In this regard, the difference between the system's state and the state estimate generated by the low-level estimator  $\hat{x}_c \in \mathbb{R}^{n_x}$ , is denoted as  $\eta = x - \hat{x}_c$ .

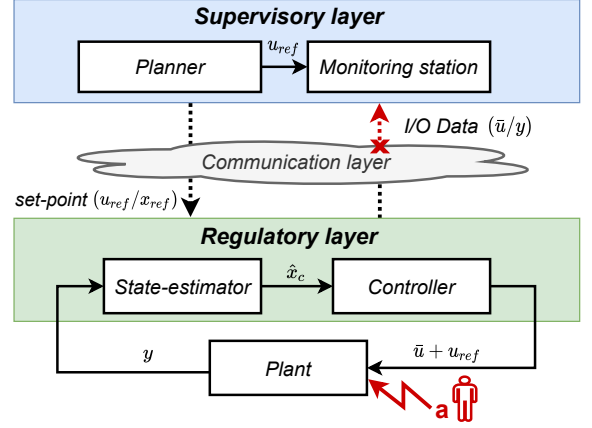


Figure 1: Overall scheme. Solid (dotted) lines represent local (remote) connections.

**Assumption 3.** The controller estimation error  $\eta$  lies within the zonotopic set

$$\eta \in \mathcal{H} = \langle c_\eta, H_\eta \rangle, \quad (7)$$

with generators matrix  $H_\eta \in \mathbb{R}^{n_x \times m_\eta}$  and center  $c_\eta \in \mathbb{R}^{n_x}$ .

Hence, the control law is given by

$$u = \bar{u} + u_{ref}, \quad (8)$$

where  $\bar{u}$  denotes the estimate-feedback action computed by the low-level controller

$$\bar{u} = -K(\hat{x}_c - x_{ref}) = -Kz - K\eta. \quad (9)$$

Accordingly, under control law (8), the tracking error dynamics is governed by the equation

$$z^+ = x^+ - x_{ref}^+ = (A - BK)z - BK\eta + E_w w, \quad (10)$$

with  $K$  designed such that  $A - BK$  is asymptotically stable. Therefore, Eq. (10) represents an asymptotically stable dynamical system subject to bounded disturbances.

#### 3.2. Supervisory layer

In the supervisory layer, an anomalies detector monitors the plant operation. For this purpose, a generic Luenberger observer is considered as follows

$$\hat{x}^+ = (A - LC)\hat{x} + Bu + Ly, \quad (11)$$

where  $\hat{x} \in \mathbb{R}^{n_x}$  represents the state estimation vector.

Let us denote the estimation error as  $e = x - \hat{x}$ . Then, from (4) and (11), we obtain

$$e^+ = x^+ - \hat{x}^+ = (A - LC)e + E_w w - LE_v v, \quad (12)$$

with  $L$  designed such that  $A - LC$  is asymptotically stable. Consequently, Eq. (12) represents an asymptotically stable dynamical system subject to bounded disturbances.

**Assumption 4.** The system is in stationary operation such that  $k \geq k^*$ , with  $k^* \in \mathbb{N}^+$  a finite sample by which the trajectories of (10) and (12) have converged into their respective mRPI sets.

In the sequel, an RPI zonotopic over-approximation of the mRPI set for system (12) will be denoted as  $\mathcal{E} = \langle c_e, H_e \rangle$ , and may be computed through Algorithm 1 in the Appendix, obtaining

$$c_e = \Lambda^{-1}(E_w c_w - L E_v c_v), \quad (13a)_{230}$$

$$H_{e,j+1} = [(A - LC)H_{e,j} \quad E_w H_w \quad -L E_v H_v], \quad (13b)$$

with  $\Lambda = I - (A - LC)$  and (13b) representing the required generators matrix recursion.

### 3.3. Anomalies detector

The presence of anomalies is monitored based on the values adopted by the following residual vector

$$r = y - C\hat{x} = C(x - \hat{x}) + E_v v = C e + E_v v. \quad (14)$$

Therefore, the healthy residual zonotopic set is given by

$$\mathcal{R}_H = C\mathcal{E} \oplus E_v \mathcal{V} = \langle c_r^h, H_r^h \rangle, \quad (15)$$

using the set operations (2), it follows

$$c_r^h = C\Lambda^{-1}(E_w c_w - L E_v c_v) + E_v c_v, \quad (16a)$$

$$H_r^h = [C H_e \quad E_v H_v]. \quad (16b)$$

Hence, the following can be established

$$\begin{cases} r \in \mathcal{R}_H & \implies \text{Healthy system,} \\ \text{otherwise} & \implies \text{Something is wrong.} \end{cases}$$

## 4. Record phase

Two different attack scenarios are considered:

- **Scenario I:** the attacker has gained access to the input/output data sent by the low-level controller to the monitoring center. Thus, the recorded data sets are  $\mathcal{Y} \equiv \{y_k : k \in \mathcal{K}_{rec}\}$  and  $\mathcal{U} \equiv \{\bar{u}_k : k \in \mathcal{K}_{rec}\}$ .
- **Scenario II:** the attacker is able to access only the output sensors data, such that the recorded data set is  $\mathcal{Y} \equiv \{y_k : k \in \mathcal{K}_{rec}\}$ . Note that this scenario models also the case where the monitoring center operates based on a set of input/output sensors installed for monitoring the plant operation and the attacker has gained access to the output sensors only. These sensors may differ from the ones used by the controller to close the low-level control loop.

Let us denote with the superscripts  $r$  and  $a$  the state of the system variables during the record and replay phases, respectively. As an example, for the state variable it follows:  $x_k^r = x_k \forall k \in \mathcal{K}_{rec}$ , while  $x_k^a = x_k \forall k \in \mathcal{K}_{rep}$ .

### 4.1. Regulatory layer

The control signal that is being injected during the record phase  $u^r = \bar{u}^r + u_{ref}^r$ , encompasses the low-level control action

$$\bar{u}^r = -Kz^r - K\eta^r, \quad (17)$$

with  $\eta^r \in \mathcal{H}$ , plus the fixed-set point reference signal  $u_{ref}^r$ .

Therefore, the tracking error dynamics is

$$z^{r+} = x^{r+} - x_{ref}^{r+} = (A - BK)z^r - BK\eta^r + E_w w^r. \quad (18)$$

### 4.2. Supervisory layer

The estimator dynamics during the record phase is

$$\hat{x}^{r+} = (A - LC)\hat{x}^r + Bu^r + Ly^r, \quad (19)$$

starting at the initial state  $\hat{x}_{k_0}^r = \hat{x}_{k_0}$ . Accordingly, the associated residual vector during the replay phase is

$$r^r = y^r - C\hat{x}^r = C(x^r - \hat{x}^r) + E_v v^r, \quad (20)$$

which, by means of Assumption 4, satisfies

$$r^r \in \mathcal{R}_H \quad \forall k \in \mathcal{K}_{rec}. \quad (21)$$

## 5. Attack phase

During the attack phase, the attacker replays back previous measurements aiming to make a physical attack conducted over the plant undetectable (cf. Section 2.3).

**Assumption 5.** The attacker knows the model of the system and is capable of compromising the state variables independently by means of an attack signal  $a(k) \in \mathbb{R}^{n_x}$  (e.g. liquid theft from different tanks of a distribution network).

Assumption 5 is in line with the attack modelling performed in other works like [3]. According to the time windows defined in Section 2.3, the attack vector satisfies

$$a(k) \begin{cases} \neq 0 & \text{if } k \in \mathcal{K}_{phy}, \\ = 0 & \text{otherwise.} \end{cases} \quad (22)$$

Following the previously presented notation, for all  $k \in \mathcal{K}_{rep}$  we have

$$x^{a+} = Ax^a + Bu^a + E_w w^a + a, \quad (23)$$

starting at the initial state  $x_{k_1}^a = x_{k_1}$ . Moreover, process disturbances/noise satisfy  $w^a \in \mathcal{W}$ ,  $v^a \in \mathcal{V}$ .

*Remark 1.* The analysis performed below is also applicable to the case in which  $a$  is injected through the input matrix, i.e. substituting  $a$  in (23) with  $B\bar{a}$  (with  $\bar{a} \in \mathbb{R}^{n_u}$ ). This case would describe cyber-attacks that modify the set-point signals sent from the supervisory to regulatory layer.

### 5.1. Regulatory layer

In the considered attacks, the malicious attacker is unable to access the dedicated network of the low-level controller, so the control loop remains healthy.

The injected signal during the replay phase is  $u^a = \bar{u}^a + u_{ref}^a$ , with

$$\bar{u}^a = -Kz^a - K\eta^a, \quad (24)$$

where  $\eta^a \in \mathcal{H}$ . Besides,  $u_{ref}^a$  denotes the reference signal that is being sent from the supervisory layer during the replay phase.

Note that the controller capability to regulate the tracking error is affected by the presence of  $a$

$$z^{a+} = x^{a+} - x_{ref}^{a+} = (A - BK)z^a + E_w w^a - BK\eta^a + a, \quad (25)$$

that is, since the controller remains healthy it will react to the physical attack  $a$ .

### 5.2. Scenario I

The first scenario considers that the recorded data sets  $\mathcal{Y}$  and  $\mathcal{U}$  are replayed back. Hence, from the supervisory layer point of view, the received signals are  $y^a = y^r$  and  $\bar{u}^a = \bar{u}^r \forall k \in \mathcal{K}_{rep}$ . Accordingly, the residual vector during the replay phase is

$$r^a = y^r - C\hat{x}^a = (y^r - C\hat{x}^r) + (C\hat{x}^r - C\hat{x}^a) = r^r + C(\hat{x}^r - \hat{x}^a), \quad (26)$$

with the state estimator evolving according to

$$\hat{x}^{a+} = (A - LC)\hat{x}^a + B\bar{u}^r + Bu_{ref}^a + Ly^r. \quad (27)$$

Let us denote by  $\bar{x} = \hat{x}^r - \hat{x}^a$  the difference in the estimation between record and replay phases. Thus, comparing (19) and (27), we obtain

$$\bar{x}^+ = \hat{x}^{r+} - \hat{x}^{a+} = (A - LC)\bar{x} + B\Delta u_{ref}, \quad (28)$$

where  $\Delta u_{ref} = u_{ref}^r - u_{ref}^a$  represents the difference in the reference signal between the record and replay phases for a fixed  $u_{ref}^r$ . Since (28) is not affected by uncertainties, by denoting  $c_{\bar{x}} = \bar{x}$  the evolution of (28) can be rewritten in zonotopic form as  $\bar{\mathcal{X}} = \langle c_{\bar{x}}, 0 \rangle$ .

Hence, taking into consideration (21), from (26) the computation of the residual set under attack is

$$\mathcal{R}_A = \mathcal{R}_H \oplus C\bar{\mathcal{X}} = \langle c_r^h, H_r^h \rangle \oplus \langle Cc_{\bar{x}}, 0 \rangle = \langle c_r^h + \delta c_r, H_r^h \rangle, \quad (29)$$

with  $\delta c_r = c_r^a - c_r^h = Cc_{\bar{x}}$  denoting the center difference.

**Definition 3.** Guaranteed attack detection is achieved if and only if  $\mathcal{R}_H \cap \mathcal{R}_A = \emptyset$  at some  $k \in \mathcal{K}_{rep}$ .

#### 5.2.1. Steady-state analysis

The main advantage of the zonotopic invariant set analysis is that it allows to derive analytic expressions regarding the separability of the residual sets in order to enforce detectability. Accordingly, below it is considered that the reference signal imposed from the supervisory layer  $u_{ref}^a$  is constant, and thus,  $\Delta u_{ref} = \text{const.}$  since  $u_{ref}^r$  is fixed.

Related to the set separation condition introduced in Definition 3, the zonotopic interpretation (see [30]) of Lemma 2.1 in [31], is formulated as

**Lemma 1.** Let  $\mathcal{Z} = \langle a_z + b_z, H_z \rangle$  and  $\mathcal{Y} = \langle a_y + b_y, H_y \rangle$ . Then,  $\mathcal{Z} \cap \mathcal{Y} = \emptyset$  if and only if  $a_y - a_z \notin \langle b_z, H_z \rangle \oplus \langle -b_y, H_y \rangle$ .

Accordingly, the following proposition regarding the output set-point imposed from the supervisory layer can be obtained.

**Proposition 1.** Guaranteed attack detection is achieved in the steady-state if the output set-point difference between record and replay phases  $\Delta y_{ref} = y_{ref}^r - y_{ref}^a$  fulfills

$$\Delta y_{ref} \notin \langle 0, M^{-1} [H_r^h \ H_r^h] \rangle, \quad (30)$$

with  $M = (I - C\Lambda^{-1}L)$ .

*Proof.* From  $\Delta u_{ref} = \text{const.}$ , and taking into consideration the asymptotically stable system (28) and that  $\delta c_r = Cc_{\bar{x}}$ , the displacement of the residual set center settles at

$$\delta c_r = c_r^a - c_r^h = C\Lambda^{-1}B\Delta u_{ref}. \quad (31)$$

By denoting as  $y_{ref}^a$  the fixed set-point generated by  $y_{ref}^a = C(I - A)^{-1}Bu_{ref}^a$ . Then, from the linearity of the reference model (6), it follows

$$\Delta y_{ref} = C(I - A)^{-1}B\Delta u_{ref}. \quad (32)$$

Besides, by taking into consideration the equality

$$\Lambda^{-1} = (I - A)^{-1} - \Lambda^{-1}LC(I - A)^{-1}, \quad (33)$$

then, using (32) and (33), (31) can be rewritten as

$$\begin{aligned} \delta c_r &= C((I - A)^{-1} - \Lambda^{-1}LC(I - A)^{-1})B\Delta u_{ref} = \\ &= (I - C\Lambda^{-1}L)\Delta y_{ref} = M\Delta y_{ref}. \end{aligned} \quad (34)$$

Therefore, by considering  $\mathcal{R}_H = \langle c_r^h, H_r^h \rangle$  and  $\mathcal{R}_A = \langle c_r^h + \delta c_r, H_r^h \rangle$ , from Lemma 1 it follows that  $\mathcal{R}_H \cap \mathcal{R}_A = \emptyset$  if and only if (30) is satisfied.  $\square$

Note that the vector  $\delta c_r$  is independent of the physical attack  $a$ , and thus its presence is masked to the supervisory layer. Besides, for  $\delta c_r = 0$ , the attack is completely undetectable in the steady-state since for this case  $\mathcal{R}_H = \mathcal{R}_A = \langle 0, H_r \rangle$ .

*Remark 2.* Based on the performed analysis, watermarking signals can be further developed by designing input sequences that take into account the transient behaviour of  $\delta c_r$  in order to enforce the guaranteed detection condition in Definition 3, while minimizing the performance loss induced in the system operation.

#### 5.3. Scenario II

For the second scenario, the received sensors data at the supervisory layer are  $y^a = y^r$ , while the controller inputs are received unaltered. Thus, the residual vector is

$$r^a = y^r - C\hat{x}^a = C(x^r - \hat{x}^a) + E_v v^r. \quad (35)$$

Denoting  $\tilde{x} = x^r - \hat{x}^a$ , its dynamics evolves according to

$$\begin{aligned} \tilde{x}^+ &= x^{r+} - \hat{x}^{a+} = (A - LC)\tilde{x} + B\Delta u_{ref} + \\ &+ B(\bar{u}^r - \bar{u}^a) + E_w w^r - LE_v v^r. \end{aligned} \quad (36)$$

Note that the evolution of (36) depends also on the dynamics of systems (18) and (25) through the control action  $\bar{u}^r$  and  $\bar{u}^a$ , respectively. Let us gather the evolution of those systems in  $q = [\tilde{x}^T \ z^r{}^T \ z^a{}^T]^T$ , such that

$$q^+ = \Theta q + \Pi d + \Sigma \Delta u_{ref} + \Phi a, \quad (37)$$

where vector  $d = [w^{rT} \ w^{aT} \ v^{rT} \ \eta^{rT} \ \eta^{aT}]^T$  encompasses the different disturbances. The augmented system matrices are

$$\Theta = \begin{bmatrix} A - LC & -BK & BK \\ 0 & A - BK & 0 \\ 0 & 0 & A - BK \end{bmatrix}, \quad \Sigma = \begin{bmatrix} B \\ 0 \\ 0 \end{bmatrix}, \quad \Pi = \begin{bmatrix} E_w & 0 & -LE_v & -BK & BK \\ E_w & 0 & 0 & -BK & 0 \\ 0 & E_w & 0 & 0 & -BK \end{bmatrix}, \quad \Phi = \begin{bmatrix} 0 \\ 0 \\ I \end{bmatrix}.$$

### 5.3.1. Steady-state analysis

An analysis similar to the one performed in Section 5.2 will be carried out. In this regard, let us consider  $u_{ref}^a$  to be constant during the replay phase, i.e.  $\Delta u_{ref} = \text{const.}$ , and let us consider the following assumption.

**Assumption 6.** The physical attack  $a$  is performed abruptly and is kept constant over the attack set  $\mathcal{K}_{phy}$ .

Consequently, at steady-state, a zonotopic over-approximation  $\mathcal{Q} = \langle c_q, H_q \rangle$  of the mRPI set for system (37) may be computed through Algorithm 1, obtaining

$$c_q = [I - \Theta]^{-1} (\Pi c_d + \Sigma \Delta u_{ref} + \Phi a), \quad H_{q,j+1} = [\Theta H_{q,j} \quad \Pi \text{diag}(H_w, H_w, H_v, H_\eta, H_\eta)],$$

where  $c_d = [c_w^T \ c_w^T \ c_v^T \ c_\eta^T \ c_\eta^T]^T$  and

$$[I - \Theta]^{-1} = \begin{bmatrix} \Lambda^{-1} & -\Lambda^{-1}BK\Gamma^{-1} & \Lambda^{-1}BK\Gamma^{-1} \\ 0 & \Gamma^{-1} & 0 \\ 0 & 0 & \Gamma^{-1} \end{bmatrix},$$

with  $\Gamma = I - (A - BK)$ .

**Proposition 2.** Guaranteed attack detection is achieved in the steady-state if the output set-point difference  $\Delta y_{ref}$  and attack vectors  $a$  satisfy

$$M \Delta y_{ref} + C \Lambda^{-1} B K \Gamma^{-1} a \notin \langle 0, [H_r^h \ H_r^a] \rangle, \quad (39)$$

with  $H_r^a = [C P H_q \ E_v H_v]$  and  $P = [I \ 0 \ 0]$ .

*Proof.* Given an attack vector  $a$  that satisfies Assumption 6 and  $\Delta u_{ref}$ , by defining the projection matrix  $P = [I \ 0 \ 0]$ , the trajectories of (36) will converge into the zonotopic set  $\tilde{\mathcal{X}} = P \mathcal{Q} = \langle c_{\tilde{x}}, H_{\tilde{x}} \rangle = \langle P c_q, P H_q \rangle$ , with

$$c_{\tilde{x}} = \Lambda^{-1} (E_w c_w - (B E_u + L E_v) c_v + B \Delta u_{ref} + B K \Gamma^{-1} a). \quad (40)$$

Hence, from (35), and taking into account (2), the residuals under attack will settle in the set

$$\mathcal{R}_A = C \tilde{\mathcal{X}} \oplus E_v \mathcal{V} = \langle c_r^a, H_r^a \rangle = \langle C c_{\tilde{x}} + E_v c_v, [C H_{\tilde{x}} \ E_v H_v] \rangle. \quad (41)$$

Accordingly, by recalling that  $c_r^h = C \Lambda^{-1} (E_w c_w - L E_v c_v) + E_v c_v$ , the center of the residual set under attack can be rewritten as a function of the healthy center as  $c_r^a = c_r^h + \delta c_r$ , where

$$\delta c_r = C \Lambda^{-1} (B K \Gamma^{-1} a + B \Delta u_{ref}). \quad (42)$$

Therefore, adapting the steps given in proof of Proposition 1, it follows that  $\mathcal{R}_H \cap \mathcal{R}_A = \emptyset$  if and only if (39) is satisfied.  $\square$

Concerning this attack scenario, the following discussion may be given regarding the attack detectability.

**Residual set size:** Note that for the attack case, the generators matrix contains additional terms with respect to the ones included in the healthy case. This is a direct consequence of the fact that the cause-effect relationship between the injected control signal and the obtained measurements during the attack is lost, i.e., the healthy control signal  $\bar{u}^a$  and the replayed output  $y^r$  take independent values during the attack. The bigger size of  $\mathcal{R}_A$  with respect to  $\mathcal{R}_H$  has two consequences: I) it is possible to detect the attack even without forcing the center displacement; II) the bigger size of the attacked set requires a bigger center displacement in order to fulfill condition (39).

**Center displacement:** Note that the attack vector  $a$  appears explicitly in the detectability condition (39). If the output set-point is maintained constant between phases  $\Delta y_{ref} = 0$ , the effect of the injected vector  $a$  is particularly critical along the directions that belong to the null space of the matrix  $C \Lambda^{-1} B K \Gamma^{-1}$ , i.e.,  $a \in \mathcal{N}(C \Lambda^{-1} B K \Gamma^{-1})$ , since these attacks would not cause a displacement of  $\delta c_r$ . In other words, a malicious attacker could carry out an unbounded attack  $a$  for which there are no detectability guarantees from the defender's point of view.

Regarding the existence of  $\mathcal{N}(C \Lambda^{-1} B K \Gamma^{-1})$ , the following proposition can be derived.

**Proposition 3.** The dimension  $d$  of  $\mathcal{N}(C \Lambda^{-1} B K \Gamma^{-1})$  is lower bounded by  $d = n_x - \text{rank}(C \Lambda^{-1} B K \Gamma^{-1}) \geq n_x - \min\{\text{rank}(C), \text{rank}(BK)\}$ .

*Proof.* The proof is based on well-known matrix rank properties. Let us denote  $X = C \Lambda^{-1}$ ,  $Y = B K \Gamma^{-1}$  such that  $\mathcal{N}(C \Lambda^{-1} B K \Gamma^{-1}) = \mathcal{N}(XY)$ . The following holds

$$\begin{aligned} \text{rank}(XY) &\leq \min\{\text{rank}(X), \text{rank}(Y)\}, \\ \text{rank}(X) &= \text{rank}(C), \\ \text{rank}(Y) &= \text{rank}(BK). \end{aligned}$$

Therefore, it follows that

$$\text{rank}(XY) \leq \min\{\text{rank}(C), \text{rank}(BK)\}.$$

Finally, by considering that the dimension of  $\mathcal{N}(XY)$  is  $d = n_x - \text{rank}(XY)$ , the proof is completed.  $\square$

Note that, given a null space  $\mathcal{N}(C \Lambda^{-1} B K \Gamma^{-1})$  of dimension  $d$ , a malicious attacker could introduce an attack  $a \in \mathcal{N}(C \Lambda^{-1} B K \Gamma^{-1})$  with  $n_x - d + 1$  components different than zero, i.e. the attacker needs to have access only to  $n_x - d + 1$  states to carry out this attack. Besides, the lower bound on  $d$  does not depend on the supervisory observer gain  $L$ . Consequently, this motivates to modify  $\Delta y_{ref}$  in order to detect these possible unbounded attacks.

**Remark 3.** Note that, since the analysis developed considers the steady-state operation of the system, the obtained

expressions are independent of the record/replay starting times.

## 6. Case Study

The considered case study is a quadruple-tank process [32], regulated using a low-level state estimated-feedback controller and supervised by means of a state estimator. False data are replayed to the anomalies detector.

Regarding the quadruple-tank system, vectors  $h = [h_1, h_2, h_3, h_4]^T$ ,  $u = [v_1, v_2]^T$  and  $y = [y_1, y_2]^T$  denote the tank levels, process inputs (voltages to the pumps) and outputs (voltages in level measurements), respectively. The system model is linearised at the minimum-phase point  $(h^*, u^*)$ , for the system parameters presented in [32].

By performing an Euler discretization with sampling time  $T_s = 1s$ , and considering the disturbance/noise input matrices  $E_w = 10^{-3}diag(5, 3, 1, 5)$ ,  $E_v = 10^{-3}diag(4, 1)$ , a discrete-time model is obtained as

$$\begin{aligned}\Delta h^+ &= A\Delta h + B\Delta u + E_w w, \\ y &= C\Delta h + E_v v,\end{aligned}$$

where  $\Delta h = h - h^*$ ,  $\Delta u = u - u^*$ . Process disturbances and sensor noise, which take random values at each sample time within the sets  $w \in \langle 0, I \rangle$  and  $v \in \langle 0, I \rangle$  are included in the simulations.

An LQR controller is designed with state and input weight matrices  $Q = 100I$  and  $R = I$ . Besides, the error for the state estimate used by the controller is constrained to  $\eta \in \langle 0, 10^{-3}I \rangle$ .

The observer in charge of the plant monitoring is also computed in an optimal way following the dual LQR design with  $Q = 100I$   $R = I$ . Setting  $\epsilon = 10^{-4}$ , a guaranteed  $\epsilon$ -approximation for the estimation error is obtained for  $l \geq 182$  (see Appendix A). By performing an interval over-approximation of the healthy residual zonotope  $\mathcal{R}_H$ , the following bounds are obtained:  $|r_1| \leq 0.0117$  and  $|r_2| \leq 0.0059$ .

### 6.1. Scenario I

The considered attack windows are:  $\mathcal{K}_{rec} = [100, 300]$  and  $\mathcal{K}_{rep} = [400, 1000]$ . Besides, an attack vector  $a = [1, 1, 1, 1]^T$  is injected during  $\mathcal{K}_{phy} = \mathcal{K}_{rep}$ . Note that the replayed data encompass the repetition of  $n_r = 3$  times the recorded data set.

The set-point differences  $\Delta y_{ref} = [\delta y_1, \delta y_2]^T$  obtained by computing an interval overapproximation of the zonotope in condition (30) are:  $|\delta y_1| \leq 1.984$ ,  $|\delta y_2| \leq 1.435$ . This means that by imposing a set-point difference in any of the outputs bigger than the computed limits  $|\delta y_1|$  and  $|\delta y_2|$ , it can be guaranteed that the residual vector will exit the healthy residual set in the steady-state.

In this regard, Figs. 2 and 3 show the system residuals and the imposed set-point for the attack described above. The dashed red lines in Fig. 2 show the computed limits of the healthy residual set. Note that during the time interval before the set-point modification (yellow background),

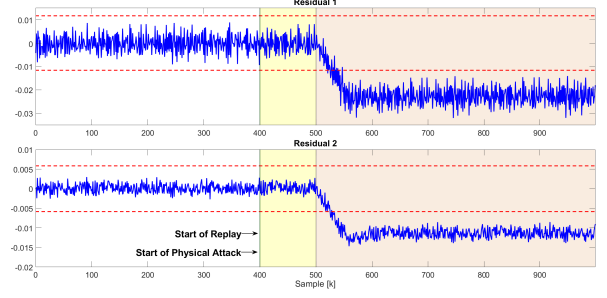


Figure 2: Residuals at the supervisory layer - Scenario I

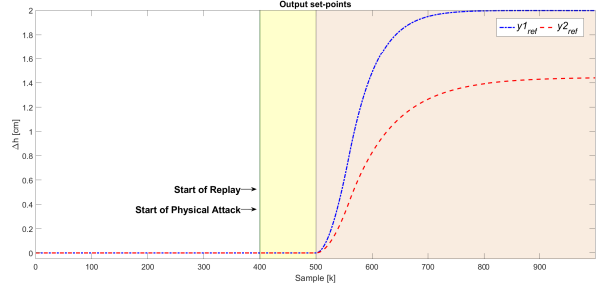


Figure 3: Set-point imposed from the supervisory layer

the attack remains completely undetectable despite the injection of vector  $a$ . Besides, it can be seen how the imposed set-point difference  $\Delta y_{ref} = [1.99, 1.44]^T$  presented in Fig. 3, enforces the system residuals to exit the healthy residual set at steady-state.

### 6.2. Scenario II

The considered attack windows in the simulation of Scenario II are:  $\mathcal{K}_{rec} = [100, 300]$  and  $\mathcal{K}_{rep} = [400, 1000]$ . For the system under study,  $\mathcal{N}(CA^{-1}BKT^{-1})$  has a dimension  $d = 2$ . Thus, the system is attacked following the direction  $a = [-0.018 \ -0.002 \ 0.017 \ 0]^T \in \mathcal{N}(CA^{-1}BKT^{-1})$  for the time interval  $\mathcal{K}_{phy} = [500, 1000]$  (note that the fourth state in vector  $a$  is set to zero).

Fig. 4 shows the effect that the injection of the attack  $a$  has on the system outputs. The real outputs (in blue) are compared to the replayed outputs (green). The attack  $a$  is introduced incipiently in the interval  $[500, 750]$  and later maintained constant. Consequently, Fig. 5 plots the residual signals generated in the supervisory layer. Since the injection of the attack  $a \in \mathcal{N}(CA^{-1}BKT^{-1})$  does not displace the attack residual center, attack detectability cannot be guaranteed unless a temporal mismatch is forced in the reference signal generated at the supervisory layer.

## 7. Conclusions

This work used zonotopic sets to develop a set-invariance analysis on the detectability of replay attacks against the supervisory layer using an observer-based detection scheme. In spite of its inherent conservativeness,



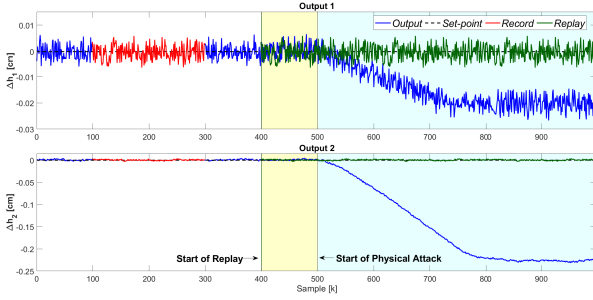


Figure 4: System outputs during the different attack phases

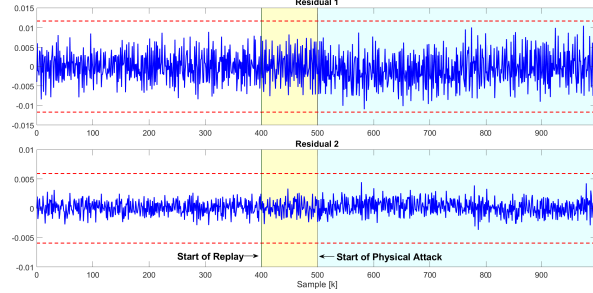


Figure 5: Residuals at the supervisory layer - Scenario II

invariant analysis is an interesting tool in the attack analysis, as it allows to derive analytical expressions regarding attack detectability. Attack detectability when an attacker is replaying directly false data has been analysed. It was shown how even in the case where the attacker is able to replay only sensor measurements, no guarantees regarding attack detectability can be given unless a temporal mismatch between record and replay phases is forced by means of a signal sent from the supervisory layer. The performed analysis serves as a basis for the future design of efficient watermarking signals that guarantee the attack detection during the transient, while minimizing the performance degradation induced to the system.

## References

- [1] H. S. Sánchez, D. Rotondo, T. Escobet, V. Puig, J. Quevedo, Bibliographical review on cyber attacks from a control oriented perspective, *Annual Reviews in Control* 48 (2019) 103–128.
- [2] A. Teixeira, I. Shames, H. Sandberg, K. H. Johansson, A secure control framework for resource-limited adversaries, *Automatica* 51 (2015) 135–148.
- [3] F. Pasqualetti, F. Dörfler, F. Bullo, Attack detection and identification in cyber-physical systems, *IEEE transactions on automatic control* 58 (11) (2013) 2715–2729.
- [4] A. A. Cárdenas, S. Amin, S. Sastry, Research challenges for the security of control systems., in: *HotSec*, 2008.
- [5] N. Hashemi, C. Murguia, J. Ruths, A comparison of stealthy sensor attacks on control systems, in: *Annual American Control Conference (ACC)*, IEEE, 2018, pp. 973–979.
- [6] Y. Liu, P. Ning, M. K. Reiter, False data injection attacks against state estimation in electric power grids, *ACM Transactions on Information and System Security (TISSEC)* 14 (1) (2011) 13.
- [7] A. Teixeira, I. Shames, H. Sandberg, K. H. Johansson, Revealing stealthy attacks in control systems, in: *50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, IEEE, 2012, pp. 1806–1813.
- [8] R. S. Smith, A decoupled feedback structure for covertly appropriating networked control systems, *IFAC Proceedings Volumes* 44 (1) (2011) 90–95.
- [9] Y. Mo, B. Sinopoli, Secure control against replay attacks, in: *47th annual Allerton conference on communication, control, and computing (Allerton)*, IEEE, 2009, pp. 911–918.
- [10] F. Miao, M. Pajic, G. J. Pappas, Stochastic game approach for replay attack detection, in: *52nd IEEE conference on decision and control*, IEEE, 2013, pp. 1854–1859.
- [11] R. M. Ferrari, A. M. Teixeira, Detection and isolation of replay attacks through sensor watermarking, *IFAC-PapersOnLine* 50 (1) (2017) 7363–7368.
- [12] C. Fang, Y. Qi, P. Cheng, W. X. Zheng, Optimal periodic watermarking schedule for replay attack detection in cyber-physical systems, *Automatica* 112 (2020) 108698.
- [13] M. Zhu, S. Martinez, On the performance analysis of resilient networked control systems under replay attacks, *IEEE Transactions on Automatic Control* 59 (3) (2013) 804–808.
- [14] M. Zhu, S. Martínez, On distributed constrained formation control in operator-vehicle adversarial networks, *Automatica* 49 (12) (2013) 3571–3582.
- [15] A. M. Mohan, N. Meskin, H. Mehrjerdi, A comprehensive review of the cyber-attacks and cyber-security on load frequency control of power systems, *Energies* 13 (15) (2020) 3860.
- [16] F. Stoican, Fault tolerant control based on set-theoretic methods., Ph.D. thesis, Supélec (2011).
- [17] C. Murguia, I. Shames, J. Ruths, D. Nešić, Security metrics and synthesis of secure control systems, *Automatica* 115 (2020) 108757.
- [18] E. Kontouras, T. Anthony, L. Dritsas, Set-theoretic detection of data corruption attacks on cyber physical power systems, *Journal of Modern Power Systems and Clean Energy* 6 (5) (2018) 872–886.
- [19] W. Lucia, B. Sinopoli, G. Franze, A set-theoretic approach for secure and resilient control of cyber-physical systems subject to false data injection attacks, in: *2016 Science of Security for Cyber-Physical Systems Workshop (SOSCYPs)*, IEEE, 2016, pp. 1–5.
- [20] G. Franzè, F. Tedesco, W. Lucia, Resilient control for cyber-physical systems subject to replay attacks, *IEEE Control Systems Letters* 3 (4) (2019) 984–989.
- [21] F. Blanchini, S. Miani, *Set-theoretic methods in control*, Springer, 2008.
- [22] V. T. H. Le, C. Stoica, T. Alamo, E. F. Camacho, D. Dumur, *Zonotopes: From guaranteed state-estimation to control*, John Wiley & Sons, 2013.
- [23] W. Kühn, Rigorously computed orbits of dynamical systems without the wrapping effect, *Computing* 61 (1) (1998) 47–67.
- [24] I. Kolmanovsky, E. G. Gilbert, Theory and computation of disturbance invariant sets for discrete-time linear systems, *Mathematical problems in engineering* 4 (1998).
- [25] F. Blanchini, Set invariance in control, *Automatica* 35 (11) (1999) 1747–1767.
- [26] G. F. Franklin, J. D. Powell, A. Emami-Naeini, J. D. Powell, *Feedback control of dynamic systems*, Vol. 4, Prentice hall Upper Saddle River, 2002.
- [27] C. Combastel, Zonotopes and Kalman observers: Gain optimality under distinct uncertainty paradigms and robust convergence, *Automatica* 55 (2015) 265–273.
- [28] T. Alamo, J. M. Bravo, E. F. Camacho, Guaranteed state estimation by zonotopes, *Automatica* 41 (6) (2005) 1035–1043.
- [29] Z. Wang, C.-C. Lim, Y. Shen, Interval observer design for uncertain discrete-time linear systems, *Systems & Control Letters* 116 (2018) 41–46.
- [30] J. K. Scott, R. Findeisen, R. D. Braatz, D. M. Raimondo, Input design for guaranteed fault diagnosis using zonotopes, *Automatica* 50 (6) (2014) 1580–1589.
- [31] D. Dobkin, J. Hershberger, D. Kirkpatrick, S. Suri, Computing the intersection-depth of polyhedra, *Algorithmica* 9 (6) (1993)



518–533.

- [32] K. H. Johansson, The quadruple-tank process: A multivariable laboratory process with an adjustable zero, *IEEE Transactions on control systems technology* 8 (3) (2000) 456–465.
- [33] S. Olaru, J. A. De Doná, M. M. Seron, F. Stoican, Positive invariant sets for fault tolerant multisensor control schemes, *International Journal of Control* 83 (12) (2010) 2622–2640.
- [34] S. V. Rakovic, E. C. Kerrigan, K. I. Kouramas, D. Q. Mayne, Invariant approximations of the minimal robust positively invariant set, *IEEE Transactions on Automatic Control* 50 (3) (2005) 406–410.
- [35] J. M. B. Caro, Control predictivo no lineal robusto basado en técnicas intervalares, Ph.D. thesis, Universidad de Sevilla (2004).
- [36] M. Althoff, O. Stursberg, M. Buss, Computing reachable sets of hybrid systems using a combination of zonotopes and polytopes, *Nonlinear analysis: hybrid systems* 4 (2) (2010) 233–249.

## Appendix A. Zonotopic $\epsilon$ -approximation of the mRPI set

Below, the computation of a zonotopic  $\epsilon$ -approximation of the mRPI set for the disturbed system (3) is discussed. This computation follows the iterative procedure used in [33] where, starting from an initial RPI, its reachable set is computed recursively, thus obtaining at each iteration a tighter RPI outer-approximation of the mRPI.

### Initial RPI set

Let us consider system (3) and the zero centered disturbance  $\delta \in \bar{\Delta} = \langle 0, H_{\Delta} \rangle$ . From Theorem 1 in [34], it follows that for a given scalar  $\alpha \in [0, 1)$  there exists a finite  $s \in \mathbb{N}^+$  that satisfies

$$A^s B \bar{\Delta} \subseteq \alpha B \bar{\Delta}. \quad (\text{A.1})$$

Besides, if (A.1) is satisfied, then the zonotope  $\langle 0, H_0 \rangle$  with

$$H_0 = (1 - \alpha)^{-1} [B H_{\Delta} \ A B H_{\Delta} \ \dots \ A^{s-1} B H_{\Delta}], \quad (\text{A.2})$$

is an RPI set for (3). Note that the evaluation of (A.1) can be formulated as a convex problem as proposed in [35].

For the specific case where matrix  $A$  in (3) has real eigenvalues, a zonotopic RPI can be obtained directly by making use of the ultimate bound analytic formula reported below.

**Theorem 1** (see [33]). Consider (3) and let  $A = V \Lambda V^{-1}$  be the Jordan decomposition of  $A$ . Then the set

$$\{x \in \mathbb{R}^n : |V^{-1}x| \leq (I - |\Lambda|)^{-1} |V^{-1}B| \bar{\delta} + \theta\}, \quad (\text{A.3})$$

is an RPI and it is attractive for the trajectories of (3), with  $\theta$  any (arbitrarily small) vector with positive elements and vector  $\bar{\delta}$  with elements  $\bar{\delta}_i = \|H_{\Delta_i}\|_1$ .

For the case in which  $A$  has real eigenvalues, the similarity transformation matrix is such that  $V \in \mathbb{R}^{n \times n}$ , and thus (A.3) is the half-space representation of a parallelepiped, which is known to be a first order zonotope. The relationship between a zero centered parallelepiped like (A.3) and its generators representation  $\langle 0, H_0 \rangle$  is formulated in [36].

### Forward propagation and stopping criterion

**Proposition 4** (see [33]). Consider (3) and denote as  $\Phi_0$  an RPI initial set for (3). Each of the set iterations:

$$\Phi_{j+1} = A\Phi_j \oplus B\Delta,$$

where  $j \in \mathbb{N}$  denotes the  $j^{\text{th}}$  element of the sequence, is an RPI approximation of the mRPI set. Moreover, as  $j$  tends to infinity, the set sequence converges to the mRPI set.

By means of the previous recursion, a certified outer  $\epsilon$ -approximation of the mRPI set  $\Omega_m$  can be obtained.

**Theorem 2** (see Theorem 3.5 in [33]). For all  $\epsilon > 0$  there exists an  $l \in \mathbb{N}^+$  such that the following RPI outer  $\epsilon$ -approximation exists:

$$\Omega_m \subset \Phi_l \subset \Omega_m \oplus \mathbb{B}_p^n(\epsilon),$$

where  $\mathbb{B}_p^n(\epsilon) = \{x \in \mathbb{R}^n : \|x\|_p \leq \epsilon\}$  and  $\|x\|_p$  is the  $p$ -norm.

From Appendix A of [33], by choosing an  $l$  such that

$$A^l \Phi_0 \subset \mathbb{B}_{\infty}^n(\epsilon/2), \quad (\text{A.4})$$

it is guaranteed that  $\Omega_m \subset \Phi_l \subset \Omega_m \oplus \mathbb{B}_{\infty}^n(\epsilon)$ . Therefore, given an  $\epsilon > 0$ , the following holds

$$\|A^l \Phi_0\|_{\infty} \leq \|A^l\|_{\infty} \|\Phi_0\|_{\infty} < \epsilon/2 \rightarrow A^l \Phi_0 \subset \mathbb{B}_{\infty}^n(\epsilon/2). \quad (\text{A.5})$$

By eigendecomposing  $A$  as  $A = T\Psi T^{-1}$ , the spectral radius of matrix  $A$  can be expressed as  $\rho(A) = \|\Psi\|_{\infty}$ . Thus,  $\|A^l\|_{\infty}$  can be bounded as

$$\|A^l\|_{\infty} = \|T\Psi^s T^{-1}\|_{\infty} \leq \|T\|_{\infty} \|T^{-1}\|_{\infty} \rho(A)^l. \quad (\text{A.6})$$

By replacing (A.6) in (A.5), and by computing  $\phi = \|\Phi_0\|_{\infty}$  and  $\kappa = \|T\|_{\infty} \|T^{-1}\|_{\infty}$ , an  $\epsilon$ -approximation to the mRPI set is guaranteed by choosing

$$l > \frac{\log(\epsilon/2) - \log(\kappa\phi)}{\log(\rho(A))}, \quad l \in \mathbb{N}^+. \quad (\text{A.7})$$

Algorithm 1 summarizes the procedure for obtaining a zonotopic  $\epsilon$ -approximation of the mRPI for the system (3).

---

### Algorithm 1 Zonotopic $\epsilon$ -approximation of the mRPI set

---

**Input:** Pair  $(A, B)$ , parameter  $\epsilon > 0$  and zonotopic representation of the disturbance set  $\Delta = \langle c_{\Delta}, H_{\Delta} \rangle$ .

**Output:** Zonotopic RPI approximation  $\mathcal{X}$  of the mRPI.

- 1: Compute  $H_0$  either using (A.2) or by means of (A.3)
  - 2: Compute the spectral radius  $\rho(A)$ ,  $\kappa$  and  $\phi = \|\Phi_0\|_{\infty}$
  - 3: Compute the minimum  $l \in \mathbb{N}^+$  such that
 
$$l > (\log(\epsilon/2) - \log(\kappa\phi)) / \log(\rho(A))$$
  - 4: For  $j = 0$  to  $j = l - 1$  propagate  $H_{j+1} = [A H_j \ B H_{\Delta}]$
  - 5: Compute the RPI set  $\mathcal{X} = \langle c_x, 0 \rangle \oplus \langle 0, H_l \rangle$  with:
 
$$c_x = (I - A)^{-1} B c_{\Delta}$$
-