




## Article

# Detection of Multiplicative False Data Injection Cyberattacks on Process Control Systems via Randomized Control Mode Switching

Shilpa Narasimhan , Matthew J. Ellis  and Nael H. El-Farra 

Department of Chemical Engineering, University of California, Davis, CA 95616, USA;  
shnarasimhan@ucdavis.edu (S.N.); mjellis@ucdavis.edu (M.J.E.)

\* Correspondence: nhelfarra@ucdavis.edu

**Abstract:** A fundamental problem at the intersection of process control and operations is the design of detection schemes monitoring a process for cyberattacks using operational data. Multiplicative false data injection (FDI) attacks modify operational data with a multiplicative factor and could be designed to be detection evading without in-depth process knowledge. In a prior work, we presented a control mode switching strategy that enhances the detection of multiplicative FDI attacks in processes operating at steady state (when process states evolve within a small neighborhood of the steady state). Control mode switching on the attack-free process at steady-state may induce transients and generate false alarms in the detection scheme. To minimize false alarms, we subsequently developed a control mode switch-scheduling condition for processes with an invertible output matrix. In the current work, we utilize a reachable set-based detection scheme and use randomized control mode switches to augment attack detection capabilities. The detection scheme eliminates potential false alarms occurring from control mode switching, even for processes with a non-invertible output matrix, while the randomized switching helps bolster the confidentiality of the switching schedule, preventing the design of a detection-evading “smart” attack. We present two simulation examples to illustrate attack detection without false alarms, and the merits of randomized switching (compared with scheduled switching) for the detection of a smart attack.

**Keywords:** cyberattack detection; multiplicative false data injection attacks; randomized control mode switching



**Citation:** Narasimhan, S.; Ellis, M.J.; El-Farra, N.H. Detection of Multiplicative False Data Injection Cyberattacks on Process Control Systems via Randomized Control Mode Switching. *Processes* **2024**, *12*, 327. <https://doi.org/10.3390/pr12020327>

Academic Editors: Scarlett Chen, Zhe Wu and Youzhi Liang

Received: 12 December 2023

Revised: 20 January 2024

Accepted: 28 January 2024

Published: 2 February 2024



**Copyright:** © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Process control systems (PCSs) are cyber-physical systems that operate chemical manufacturing processes and may utilize networked communication to integrate the cyber components (e.g., controllers and human machine interfaces) with the physical components (e.g., sensors and actuators). Cyberattacks on PCSs involve malicious incursions into the PCS network with the objective of modifying the process operational data transferred over the communication channels linking the sensors to the controller and/or the communication channels linking the controller to the actuators. Cyberattacks have increased in frequency over the past decade [1]. A loss of control in critical infrastructure, like chemical manufacturing processes, due to cyberattacks may have disastrous consequences, and therefore, the research focused on the enhancement of PCS cybersecurity has become increasingly important [2].

Due to the interconnected nature of a cyber-physical system like a PCS, a comprehensive solution to enhancing the cybersecurity of PCSs may involve adopting a multi-faceted approach that involves both information technology-based and operational technology-based approaches. The approaches proposed in the literature broadly include those that reinforce information technology-based infrastructure (e.g., firewalls [3] and network design-based approaches [4,5]), approaches that involve cybersecure architecture design

(e.g., [6]), approaches for process equipment design to mitigate adverse impacts of a cyberattack (e.g., [7,8]), and the design of operational technology (OT)-based approaches for detecting, identifying, and successfully recovering from an ongoing cyberattack [6,9–20].

The design of approaches to enable the detection of an ongoing cyberattack in particular has received significant attention. Several approaches have been proposed in the literature. Examples include detection schemes that monitor the values of a residual (defined as a vector that is the difference between a measured variable and its estimated value) and detect an attack if the residual breaches a certain threshold. Other examples include set membership-based detection schemes that monitor a process by confirming that the residual values are contained within a set that describes all possible attack-free values of the residual [14,16]. Machine learning-based approaches have also been used as methods for cyberattack detection [6,17]. Other examples include residual-based detection strategies that are designed to detect an attack if the attack causes the process states to leave a safe operating region [10–12]. Attack detection approaches that utilize random authentication signals to the manipulated input to add a perturbation to the process and reveal the presence of an attack [21–25] have also been proposed. The approach chosen for attack detection may depend on the type of attack to be detected. False data injection (FDI) attacks aim to compromise the integrity of the PCS by maliciously modifying the data over the communication channels [26,27]. Multiplicative FDI attacks alter the value of process variables over the PCS network using a multiplicative factor. These attacks can be designed to evade detection without requiring extensive process knowledge, and the design of approaches for their detection has received some attention in the literature [13,15,16,23–25].

In this work, we use a reachable set-based detection scheme [14] that tracks the evolution of some monitoring variable with respect to the evolution of the reachable sets of the attack-free process to detect an attack. The reachable sets take account of the attack-free process behavior during transient operation, and therefore the detection scheme utilizing the reachable sets generates no alarms when monitoring the process under control mode switching in the absence of attacks. The reachable set-based detection scheme makes no assumptions on the structure of the output matrix, and generates no alarms when monitoring processes with non-invertible output matrices under control mode switching in the absence of attacks. To detect a range of multiplicative FDI attacks, we propose a randomized switching-enabled cyberattack detection method under which the control mode is switched at randomly chosen switching instances. Randomization to enhance the cybersecurity of a control system has received some attention in the literature [22,23,28,29]. However, randomization to enhance the cybersecurity of a detection scheme has not been explored. Under our proposed randomized switching-based detection method, the control mode switching instances are randomly chosen to confound an attacker aiming to learn the switching schedule. Without complete knowledge of the switching schedule, an attacker may not be able to design a “smart” attack that is capable of evading detection, thereby enhancing the cybersecurity of the detection method.

The remainder of the paper is organized as follows: in Section 2, we present some preliminaries and discuss the class of processes and the attack model. Following this, in Section 3, we characterize the interdependence between the control parameters, the closed-loop stability of the process, and the detectability of a multiplicative attack with respect to the reachable set-based detection scheme. We propose a control mode switching-enabled method for the detection of multiplicative FDI attacks on a process under transient operation. We then present theoretical results to demonstrate that the randomized switching-enabled attack detection method guarantees a zero false alarm rate from multiple successive control mode switches on the process under transient operation without requiring that the output matrix be invertible. We discuss some implementation considerations and present algorithms for two variations of the switching-enabled attack detection method. In Section 4, we apply the detection method over simulations using two illustrative examples and demonstrate that the method guarantees attack detection while generating no false alarms when monitoring the process in the absence of attacks. Finally, using simulations

of the second illustrative example, where we apply the detection method to a chemical process operated at its steady state, we demonstrate that for the detection of a smart attack, the randomized switching-enabled attack detection method may be preferable to a method using scheduled control mode switches.

## 2. Preliminaries

### 2.1. Notation

$\mathbb{Z}^+$  is the set of non-negative integers.  $\mathbb{R}^n$  is the  $n$ -dimensional Euclidean space. Given a vector  $x = [x_1 \ x_2 \ \dots \ x_n]^T \in \mathbb{R}^n$ ,  $\|x\| = \sqrt{\sum_{i=1}^n (x_i)^2}$  is its Euclidean norm and  $\|x\|_\infty = \max_{i=1}^n |x_i|$  is its  $\infty$ -norm. For a square matrix  $A \in \mathbb{R}^{n \times n}$ , its spectral radius is defined as  $\rho(A) = \max\{|\lambda_1|, |\lambda_2|, \dots, |\lambda_n|\}$ , where  $\lambda_i$  is the  $i^{\text{th}}$  eigenvalue of the matrix  $A$ . The Minkowski sum of two sets  $\mathcal{X} \subset \mathbb{R}^n$  and  $\mathcal{Y} \subset \mathbb{R}^n$  is defined as  $\mathcal{X} \oplus \mathcal{Y} = \{x + y \mid x \in \mathcal{X}, y \in \mathcal{Y}\}$ . Given a matrix  $A \in \mathbb{R}^{m \times n}$  and a set  $\mathcal{X} \subset \mathbb{R}^n$ ,  $A\mathcal{X} = \{Ax \mid x \in \mathcal{X}\}$  is the linear map of the set  $\mathcal{X}$ . Given a matrix  $A \in \mathbb{R}^{m \times n}$  and a set  $\mathcal{X} \subseteq \mathbb{R}^n$ ,  $\bigoplus_{i=0}^t A^i \mathcal{X}$  represents the series  $\mathcal{X} \oplus A\mathcal{X} \oplus \dots \oplus A^t \mathcal{X}$ . Given a vector  $x \in \mathbb{R}^n$ , an  $n$ -dimensional polytope is a bounded region in the Euclidean space that satisfies the linear matrix inequalities  $Ax \leq b$ , where  $A \in \mathbb{R}^{m \times n}$  is a matrix and  $b \in \mathbb{R}^{m \times 1}$  is a vector. A zonotope is a convex polytope that is symmetric about its center and may be formally defined as the Minkowski sum of a finite set of line segments [30].

### 2.2. Class of Processes

In this work, we consider processes that are modeled by discrete-time linear time-invariant systems of the form:

$$x_{t+1} = A^x x_t + B^u u_t + B^w w_t \quad (1a)$$

$$y_t = C^x x_t + v_t \quad (1b)$$

where  $x_t \in \mathbb{R}^{n_x}$  and  $u_t \in \mathbb{R}^{n_u}$  for all  $t \in \mathbb{Z}^+$  are the vectors representing the process states and manipulated inputs, respectively, and the measured output is  $y_t \in \mathbb{R}^{n_y}$  with  $n_y \leq n_x$ . The process is subject to bounded process disturbances and measurement noise, where  $w_t \in \mathcal{W} \subset \mathbb{R}^{n_x}$  and  $v_t \in \mathcal{V} \subset \mathbb{R}^{n_y}$  are the vectors representing the bounded process disturbances and measurement noise. The compact sets  $\mathcal{W}$  and  $\mathcal{V}$  are assumed to be known polytopes.  $A^x \in \mathbb{R}^{n_x \times n_x}$ ,  $B^u \in \mathbb{R}^{n_x \times n_u}$ ,  $B^w \in \mathbb{R}^{n_x \times n_w}$ , and  $C^x \in \mathbb{R}^{n_y \times n_x}$  are matrices.

To estimate the process states, a Luenberger observer is utilized:

$$\hat{x}_{t+1} = A^x \hat{x}_t + B^u u_t + L(y_t - \hat{y}_t) \quad (2a)$$

$$\hat{y}_t = C^x \hat{x}_t \quad (2b)$$

where  $\hat{x}_t \in \mathbb{R}^{n_x}$  and  $\hat{y}_t \in \mathbb{R}^{n_y}$  for all  $t \in \mathbb{Z}^+$  are the estimates of the process states and the measured outputs generated by the observer. The observer gain  $L \in \mathbb{R}^{n_x \times n_y}$  is chosen such that all eigenvalues of the matrix  $A - LC$  are strictly within the unit circle. The estimates of the process states generated by the observer are utilized to compute the control action as follows:

$$u_t = -K\hat{x}_t \quad (3)$$

where  $K \in \mathbb{R}^{n_x \times n_u}$  is the gain of the feedback controller, which is chosen to ensure that all eigenvalues of the matrix  $A - BK$  are contained strictly inside the unit circle.

The discrepancy between the process states and their estimates generated by the observer is defined as the estimation error ( $e := x - \hat{x}$ ) and has the dynamics:

$$e_{t+1} = (A^x - LC^x)e_t + B^w w_t - Lv_t \quad (4)$$

The dynamics of the overall closed-loop process are collectively described by the evolution of the process states and the estimation errors, which we analyze using an

augmented state  $\tilde{\zeta} = [x^T \ e^T]^T$ , which consists of the process states and estimation errors. The dynamics of the closed-loop augmented state are governed by:

$$\tilde{\zeta}_{t+1} = \underbrace{\begin{bmatrix} A^x - B^u K & B^u K \\ 0 & A^x - LC^x \end{bmatrix}}_{A^{\tilde{\zeta}}(K,L)} \tilde{\zeta}_t + \underbrace{\begin{bmatrix} B^w & 0 \\ B^w & -L \end{bmatrix}}_{B^d(L)} d_t \quad (5)$$

where  $d_t := [w_t^T \ v_t^T]^T$  is an augmented disturbance vector, which consists of the process disturbances and measurement noise and is bounded within the compact set  $\mathcal{D} := \{[w^T \ v^T]^T \mid w \in \mathcal{W}, v \in \mathcal{V}\}$ . Without loss of generality, the origin is assumed to be the steady state of the process.

In the remainder of the paper, the term “attack-free closed-loop process with  $(K, L)$  and a set of initial states  $\mathcal{R}_0^{\tilde{\zeta}}$ ” is used to refer to the process with augmented state dynamics, as in Equation (5). In this work, we consider process transients during which the states of the closed-loop process in Equation (5) evolve from a compact set of initial states ( $\tilde{\zeta}_0 \in \mathcal{R}_0^{\tilde{\zeta}} \subset \mathbb{R}^{2n_x}$ ). For the closed-loop process with  $(K, L)$ , its  $t$ -step reachable set is defined as the set of all states that can be reached in  $t$  time steps from the set of initial states  $\mathcal{R}_0^{\tilde{\zeta}}$  and under all admissible disturbances [31]. The  $t$ -step reachable set of the closed-loop process may be computed using the following relationship:

$$\mathcal{R}_t^{\tilde{\zeta}}(K, L) = A^{\tilde{\zeta}}(K, L) \mathcal{R}_{t-1}^{\tilde{\zeta}}(K, L) \oplus B^d(L) \mathcal{D}, \text{ for } t > 0 \quad (6)$$

From Equation (6), the  $t$ -step reachable set is dependent on the set of initial states  $\mathcal{R}_0^{\tilde{\zeta}}$ , the control parameters  $(K, L)$ , the time step  $t \in \mathbb{Z}^+$ , and the set of bounded disturbances  $\mathcal{D}$ . For a concise representation, in this work, we highlight the dependence of the  $t$ -step reachable set on the control parameters only. When the closed-loop process with  $(K, L)$  is stable with  $\rho(A^{\tilde{\zeta}}(K, L)) < 1$ , its augmented state is ultimately bounded within the minimum invariant set of the process, which is the limit set of all trajectories of the process [32]. The minimum invariant set, which contains all possible values of augmented state of the closed-loop process under steady-state operation can be expressed as [33]:

$$\mathcal{R}_{\infty}^{\tilde{\zeta}}(K, L) = \bigoplus_{i=0}^{\infty} A^{\tilde{\zeta}}(K, L)^i B^d(L) \mathcal{D} \quad (7)$$

### 2.3. Class of Multiplicative False Data Injection Attacks

In this work, we consider the detection of multiplicative false data injection (FDI) attacks that simultaneously alter the data communicated over the sensor–controller and controller–actuator links of the PCS network. The falsified data are represented as follows:

$$y_t^a = \Lambda^y y_t \quad (8a)$$

$$u_t^a = \Lambda^u u_t \quad (8b)$$

where  $\beta_t^a \in \mathbb{R}^{n_\beta}$  is the value of  $\beta \in \mathbb{R}^{n_\beta}$  that is modified by the attack,  $\Lambda^\beta \neq I$  is the pre-multiplicative factor that alters the value of  $\beta$  in the presence of an attack. If  $\beta = y$ , the attack alters the value of the measured output communicated to the controller, and if  $\beta = u$ , the attack alters the value of the manipulated input communicated to the control actuators. In the presence of an attack,  $u^a$  is the implemented control action that is used by the observer in Equation (2a) to generate estimates of process states. An attack alters the closed-loop process dynamics in Equation (5) as follows:

$$\tilde{\zeta}_{t+1} = \underbrace{\begin{bmatrix} A^x - B^u \Lambda^u K & B^u \Lambda^u K \\ L(I - \Lambda^y) C^x & A^x - LC^x \end{bmatrix}}_{A^{\tilde{\zeta}^a}(K,L)} \tilde{\zeta}_t + \underbrace{\begin{bmatrix} B^w & 0 \\ B^w & -L \Lambda^y \end{bmatrix}}_{B^{d_a}(L)} d_t \quad (9)$$

In the remainder of the paper, the term “attacked closed-loop process with  $(K, L)$  and a set of initial states  $\mathcal{R}_0^\xi$ ” is used to refer to the process under an attack with augmented state dynamics, as in Equation (9). The  $t$ -step reachable set of the augmented state for the process under an attack with  $(K, L)$  and a set of initial states  $\mathcal{R}_0^\xi$  may be computed per the following relationship [31]:

$$\mathcal{R}_t^{\xi_a}(K, L) = A^{\xi_a}(K, L)\mathcal{R}_{t-1}^{\xi_a}(K, L) \oplus B^{d_a}(L)\mathcal{D}, \text{ for } t > 0 \quad (10)$$

where the subscript “ $a$ ” to the variable (e.g.,  $\mathcal{R}_t^{\xi_a}(K, L)$ ,  $A^{\xi_a}$ , and  $B^{d_a}$ ) indicates the presence of an attack. The minimum invariant set of the attacked closed-loop process with  $(K, L)$  stable in the sense that  $\rho(A^{\xi_a}(K, L)) < 1$ , may be expressed (similar to Equation (7)) as:

$$\mathcal{R}_\infty^{\xi_a}(K, L) = \bigoplus_{i=0}^{\infty} A^{\xi_a}(K, L)^i B^{d_a}(L)\mathcal{D} \quad (11)$$

#### 2.4. Class of Detection Schemes Utilizing Reachable Sets

In the presence of an attack on the process, the augmented state values may deviate from their expected values for the attack-free process, meaning that an anomaly in the augmented state value, if detected, may indicate the presence of an attack. However, the augmented state cannot be measured, and standard anomaly detection schemes monitor a process based on values of a monitoring variable that are a function of the augmented state. Commonly used monitoring variables include the measured output  $y = [C^x \ 0]\xi + [0 \ I]d$ , the estimated output  $\hat{y} = [C^x - C^x]\xi$ , or the residual, which may be defined as  $r := y - \hat{y} = [0 \ C^x]\xi + [0 \ I]d$ . The ability of a detection scheme to detect an attack may vary with the monitoring variable used. For example, detection of some attacks may be possible when the detection scheme uses the residual to monitor the process; however, these attacks may go undetected when the detection scheme uses the measured output. This realization motivated the formulation of a generalized monitoring variable in [13,14] which may be expressed in terms of the measured and estimated outputs as:

$$\eta_t = H^y y_t + H^{\hat{y}} \hat{y}_t \quad (12)$$

where  $H^y$  and  $H^{\hat{y}}$  are design parameters that dictate the choice of the monitoring variable. The monitoring variable  $\eta$  is a design parameter for a detection scheme that takes the form of Equation (15), in the sense that each choice of the monitoring variable gives a different reachable set-based detection scheme. For example, if  $H^y = [C^x \ 0]$  and  $H^{\hat{y}} = [0 \ I]$ , then  $\eta = y$ , and if  $H^y = [0 \ C^x]$  and  $H^{\hat{y}} = [0 \ I]$ , then  $\eta = r$ .

In view of the definitions of the augmented state and disturbance vectors, the monitoring variable may be expressed as a linear combination of the augmented state and the disturbance vector as follows:

$$\eta_t = \underbrace{[(H^y + H^{\hat{y}})C^x - H^{\hat{y}}C^x]}_{:=C^\eta} \xi_t + \underbrace{[0 \ H^y]}_{:=D^\eta} d_t \quad (13)$$

For the closed-loop process with  $(K, L)$  and a set of initial states  $\mathcal{R}_0^\xi$ , all possible values of the monitoring variable are bounded within the  $t$ -step reachable sets for all  $t \in \mathbb{Z}^+$ . From Equation (13), the  $t$ -step reachable set of the monitoring variable depends on the control parameters  $(K, L)$ , the set of initial states  $\mathcal{R}_0^\xi$ , and the set  $\mathcal{D}$  and may be computed per the following relationship:

$$\mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K, L)) = C^\eta \mathcal{R}_t^\xi(K, L) \oplus D^\eta \mathcal{D} \quad (14)$$

In this work, we use a class of detection schemes that track the evolution of the monitoring variable at each time step with respect to the reachable sets of the attack-free process at that time step [14]. The detection logic is given by:

$$\phi_t(\eta_t) = \begin{cases} 0, & \eta_t \in \mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K, L)) \\ 1, & \eta_t \notin \mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K, L)) \end{cases} \quad (15)$$

where  $\phi_t(\eta_t)$  is the output of the detection scheme. A value of  $\phi_t(\eta_t) = 1$  at time step  $t \in \mathbb{Z}^+$  indicates that the detection scheme generates an alarm because an attack has been detected, while a value of  $\phi_t(\eta_t) = 0$  indicates that no alarm is raised because an attack has not been detected. For a concise presentation of the results, we will refer to the detection scheme described by Equation (15) as the reachable set-based detection scheme. The reachable set-based detection scheme guarantees a zero false alarm rate during transient process operation because it accounts for the evolution of the monitoring variable of the attack-free process.

In [14], we presented a method for the classification of attacks on the basis of their detectability, which we defined as the ability of the detection scheme in Equation (15) to detect an attack. Here, we present an overview of the detectability-based attack classification method. If an attack on the closed-loop process with  $(K, L)$  and the set of initial states  $\mathcal{R}_0^\xi$  is detected by the detection scheme in finite time, the attack is called a detectable attack. Attacks that cannot be detected by the detection scheme in finite time are called undetectable attacks. Finally, a potentially detectable attack is any attack that cannot be classified as either detectable or undetectable.

For the process under attack, the reachable sets for the monitoring variable are given by:

$$\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K, L)) = C^{\eta_a} \mathcal{R}_t^{\xi_a}(K, L) \oplus D^{\eta_a} \mathcal{D} \quad (16a)$$

$$\eta_t^a = \underbrace{[(H^y \Lambda^y + H^{\hat{y}}) C^x - H^{\hat{y}} C^x] \xi_t}_{:= C^{\eta_a}} + \underbrace{[0 \ H^y \Lambda^y] d_t}_{:= D^{\eta_a}} \quad (16b)$$

It should be noted that, for the stable closed-loop system (both in the presence and in the absence of an attack), the monitoring variable is ultimately bounded within the minimum invariant set. The minimum invariant set of the monitoring variable is computed based on the minimum invariant set of the augmented state and the disturbance set, using the following relationship:

$$\begin{aligned} \mathcal{R}_\infty^\eta(K, L) &= C^\eta \mathcal{R}_\infty^\xi(K, L) \oplus D^\eta \mathcal{D} \\ \mathcal{R}_\infty^{\eta_a}(K, L) &= C^{\eta_a} \mathcal{R}_\infty^{\xi_a}(K, L) \oplus D^{\eta_a} \mathcal{D} \end{aligned} \quad (17)$$

where  $\mathcal{R}_\infty^\eta(K, L)$  and  $\mathcal{R}_\infty^{\eta_a}(K, L)$  are the minimum invariant sets of the monitoring variable for the attack-free and attacked processes, respectively. In the remainder of the paper, to distinguish these minimum invariant sets from those associated with the augmented state, the minimum invariant sets for the monitoring variable are referred to as the terminal sets.

From Equations (14) and (16a), it can be seen that the detectability of an attack (with respect to the detection scheme in Equation (15)) is dependent on how the reachable sets of the monitoring variable for the process under attack evolve in relation to the reachable sets of the monitoring variable for the process in the absence of an attack. Given an attack and the set of initial states, the reachable sets for the attack-free and the attacked processes may be computed offline, and attacks may be classified on the basis of their detectability with respect to the reachable set-based detection schemes by checking for certain conditions [14]. If at any time step  $t \in \mathbb{Z}^+$ , the intersection between the reachable set for the attacked process and the reachable set for the attack-free processes is empty ( $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K, L)) \cap \mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K, L)) = \emptyset$ ), then the attack is detectable. An attack is unde-



tectable if, for all  $t \in \mathbb{Z}^+$ , the reachable set for the attacked process is contained within the reachable set for the attack-free process, i.e.,  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K, L)) \subseteq \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K, L))$ . An attack is potentially detectable if for all  $t \in \mathbb{Z}^+$ , the reachable sets for the attacked process intersect with, but not necessarily contained within, the reachable sets for the attack-free process, i.e.,  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K, L)) \cap \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K, L)) \neq \emptyset$  for all  $t \in \mathbb{Z}^+$  and  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K, L)) \not\subseteq \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K, L))$  for some  $t \in \mathbb{Z}^+$ .

### 3. Randomized Control Mode Switching for Detection of Cyberattacks

In this section, we present a reachable set-based detection scheme that employs randomized switching between different control modes to facilitate attack detection. To provide context for the development of the proposed scheme, we begin by reviewing our previous work on control mode switching-enabled attack detection designed for processes under steady-state operation. Following this, we present theoretical results that characterize the interdependence between the control parameter selection, closed-loop stability under attack, and the detectability of an attack with respect to the detection scheme in Equation (15). We present an algorithm for a randomized control mode switching-enabled attack detection method designed for processes during a transient operation to enable attack detection with zero false alarms. Finally, we propose a modification to the control mode switching algorithm that enables attack detection with zero false alarms when implemented on processes under steady-state operation.

#### 3.1. Control Mode Switching for Cyberattack Detection

Multiplicative FDI attacks alter the stability properties of a closed-loop process by modifying the eigenvalues of the matrix  $A^{\xi_a}(K, L)$  in Equation (9). As a result, the detectability of such attacks may be influenced by the stability of the closed-loop process under an attack. In prior works [13,15], we presented a control mode switching-enabled attack detection method for the detection of multiplicative sensor–controller link FDI attacks. We considered processes under steady-state operation, when the augmented state and the monitoring variable are bounded within the minimum invariant set and the terminal set, respectively. The switching-enabled attack detection method enabled attack detection by exploiting the interdependence between the control parameter selection, the stability of the closed-loop process under an attack, and the detectability of an attack with respect to the terminal set-based detection scheme. Dynamics of the attack-free process may be excited due to a control mode switch and may cause a brief transient operation of the process during which the monitoring variable may not be bounded within the terminal set. The terminal set-based detection scheme may therefore generate a false alarm during transient operation, as it does not distinguish between anomalies in values of the monitoring variable during an attack and during transient operation. For processes with an invertible output matrix, we presented a condition that may be checked to schedule control mode switches at time steps when the switch does not cause a transient operation in the attack-free process [13]. However, satisfaction of the presented condition could not be guaranteed. To enable attack detection, a control mode switch could be implemented at a time when the condition is not satisfied, potentially leading to a transient operation, which could trigger false alarms in the detection scheme. Therefore, the proposed method did not eliminate false alarms.

In this work, we utilize the reachable set-based detection scheme in Equation (15) to monitor a process during transient operation to eliminate false alarms even when the output matrix is not invertible. From Equation (10), the reachable sets for the attacked process are influenced by the control parameters, and by extension, the control parameters influence attack detectability with respect to the reachable set-based detection scheme. If the attacked closed-loop process is unstable, its reachable sets may evolve differently from the reachable sets associated with the attack-free process, thereby influencing attack detectability. Proposition 1 below characterizes the interdependence between the stability of the closed-loop process under an attack and the detectability of an attack with respect to the detection scheme in Equation (15).

**Proposition 1.** Consider the attacked closed-loop process with  $(K, L)$  and a set of initial states  $\mathcal{R}_0^\xi$ . Let an attack destabilize the process in the sense  $\|\xi_t\| \rightarrow \infty$  as  $t \rightarrow \infty$ . If the matrix pair  $(A^{\xi_a}(K, L), C^{\eta_a})$  is observable, then the attack is detectable with respect to the reachable set-based detection scheme in Equation (15).

**Proof.** For an attack-free closed-loop process with  $(K, L)$  and a set of initial states  $\mathcal{R}_0^\xi$ , its  $t$ -step reachable set is a compact set. This follows from Equation (16a), and the assumptions that the disturbances are bounded within a compact set  $(\mathcal{D})$  and the set of initial states  $(\mathcal{R}_0^\xi)$  is a compact set. Let the attack cause the closed-loop process to be unstable in the sense that  $\|\xi_t\| \rightarrow \infty$  as  $t \rightarrow \infty$ . From Theorem 1 in [15], if the matrix pair  $(A^{\xi_a}(K, L), C^{\eta_a})$  is observable, then the 2-norm of the monitoring variable of the process also grows unbounded with time, i.e.,  $\|\eta_t\| \rightarrow \infty$  as  $t \rightarrow \infty$ . As a result, there exists a finite time step  $t' > 0$  at which the monitoring variable is not bounded within the  $t'$ -step reachable set of the attack-free process  $(\eta_{t'} \notin \mathcal{R}_{t'}^\eta(\mathcal{R}_{t'}^\xi(K, L)))$ , and the detection scheme in Equation (12) detects the attack at the time step  $t' > 0$  with an output of 1. Therefore, the attack is detectable with respect to the reachable set-based detection scheme in Equation (15).  $\square$

Proposition 1 establishes a sufficient condition for an attack to be detectable by the reachable set-based detection scheme. To enable attack detection, it may be preferable to operate the closed-loop process under an “attack-sensitive” mode, in which the control parameters are chosen such that an attack destabilizes the process. However, as noted in [13,15], prolonged operation of the process under the attack-sensitive control mode may be undesirable because a tradeoff between attack detection and closed-loop performance may exist. To manage the tradeoff, we consider extended operation of the process under a “nominal” control mode, for which the control parameters  $(K^N, L^N)$  are chosen to meet closed-loop performance considerations. For enabling attack detection, we consider control mode switching to operate the process under the attack-sensitive control mode, for which the control parameters  $(K^A, L^A)$  are chosen so that a range of attacks are detectable (per Proposition 1). We choose attack-sensitive control parameters so that an attack in the range considered causes the closed-loop process operated under the attack-sensitive control mode to be unstable in the sense that  $\rho(A^{\xi_a}(K^A, L^A)) > 1$  and the matrix pair  $(A^{\xi_a}(K^A, L^A), C^{\eta_a})$  is observable.

Control parameters influence the reachable sets of the attack-free process (from Equation (6)). Under the switching-enabled detection method, the control parameters may vary with the time step. Therefore, to monitor the switched system, we modify the detection scheme in Equation (15) as follows:

$$\phi_t(\eta_t) = \begin{cases} 0, & \eta_t \in \mathcal{R}_t^\eta(\mathcal{R}_t^\xi) \\ 1, & \eta_t \notin \mathcal{R}_t^\eta(\mathcal{R}_t^\xi) \end{cases} \quad (18)$$

where  $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi)$  is the reachable set of the attack-free process that is computed per the following recurrence relation:

$$\mathcal{R}_t^\xi = A^\xi(K_{t-1}, L_{t-1})\mathcal{R}_{t-1}^\xi \oplus B^d(L_{t-1})\mathcal{D}, \text{ for } t > 0 \quad (19)$$

where  $(K_{t-1}, L_{t-1})$  are the gains at time  $t - 1$ ,  $(K_t, L_t) = (K^N, L^N)$  if the control system is operated under the nominal mode, and  $(K_t, L_t) = (K^A, L^A)$  if the control system is operated under the attack-sensitive mode. The dependence of the reachable sets for the switched system on the control parameters is dropped to simplify notation, i.e., the reachable sets for the augmented state and the monitoring variable at any time step  $t \in \mathbb{Z}^+$  is denoted simply as  $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi)$  and  $\mathcal{R}_t^\xi$ . We initialize the reachable sets for the attack-free process at the set of initial states  $(\mathcal{R}_t^\xi = \mathcal{R}_0^\xi \text{ at } t = 0)$  and assume that initially the nominal control mode is the active mode, so that the first switch occurs from the nominal mode to the attack-sensitive mode. For attack detection, randomly choosing the control mode switching instance  $(t_s)$



may be preferable to scheduling the control mode switch at pre-determined time steps. This is because randomly choosing the switching instances helps preserve the confidentiality of the switching schedule and prevent an attacker from designing a “smart” attack that evades detection. The reachable set-based detection scheme in Equation (18) accounts for the evolution of the attack-free monitoring variable under each control mode and guarantees the generation of no false alarms in the presence of any control mode switching implemented at any randomly chosen switching instance ( $t_s \in \mathbb{Z}^+$ ). This is proven in the following proposition, which considers multiple switches between the nominal control mode and the attack-sensitive control mode implemented on the attack-free process at randomly chosen time steps  $t_{s_i} \in \mathbb{Z}^+$ , where  $i \in \{1, 2, 3, \dots\}$  such that  $t_{s_{i+1}} > t_{s_i}$ . The odd values of  $i$  represent the time steps at which a switch from the nominal control mode to the attack-sensitive control mode is implemented, while even values of  $i$  represent the time steps when a switch back from the attack-sensitive mode to the nominal mode occurs.

**Proposition 2.** *Consider the attack-free closed-loop process operated under the nominal control mode with a set of initial states  $\mathcal{R}_0^\xi$ , which is monitored by the reachable set-based detection scheme in Equation (18). Let multiple switches between the nominal and the attack-sensitive control modes be implemented on the process. Let all switching instances be randomly chosen time steps  $t_{s_i} \in \mathbb{Z}^+$ , where  $i \in \{1, 2, 3, \dots\}$  such that  $t_{s_{i+1}} > t_{s_i}$ . The reachable set-based detection scheme generates no alarms for all  $t \in \mathbb{Z}^+$ .*

**Proof.** Consider the attack-free closed-loop process operated under the nominal control mode with the set of initial states  $\mathcal{R}_0^\xi$ . Before a switch to the attack-sensitive control mode occurs, the reachable set-based detection scheme in Equation (18) accounts for all values of the attack-free process and generates no false alarms. Consider the switch to the attack-sensitive control mode at the randomly chosen time step  $t_{s_1} \in \mathbb{Z}^+$ , followed by another switch back to the nominal control mode at the randomly chosen time step  $t_{s_2} > t_{s_1}$ . From Equation (18), the reachable set-based detection scheme switches to monitoring the process based on the attack-free reachable sets computed for the attack-free process operated under the attack-sensitive mode, i.e.,  $(K_t, L_t) = (K^A, L^A)$  for all  $t \in [t_{s_1}, t_{s_2})$ . For the attack-free process operated under the attack-sensitive control mode, its reachable set at each time step after the switch contains all values of the monitoring variable ( $\eta_t \in \mathcal{R}_t^\eta(\mathcal{R}_i^\xi)$ ), leading to an output of 0 for  $t \in [t_{s_1}, t_{s_2})$ . Similarly, after switching back to the nominal control mode at time step  $t_{s_2}$ , the reachable set-based detection scheme generates no false alarms. Similarly, for all subsequent switches between the nominal and the attack-sensitive control modes implemented on the process at randomly chosen time steps, it can be demonstrated that the detection scheme in Equation (18) generates no alarms.  $\square$

From Proposition 2, the reachable set-based detection scheme guarantees a zero false alarm rate under multiple successive switches between the nominal and attack-sensitive control modes implemented on the attack-free process. The absence of an attack on the process is a sufficient condition to be satisfied for zero false alarms, meaning that a lack of alarms may not be indicative of attack-free process behavior. However, it follows from Proposition 2 that if the reachable set-based detection scheme generates an alarm at any time step  $t_d \in \mathbb{Z}^+$ , then the alarm is only due to an attack on the process.

In general, it is not known if an attack is occurring on the process. To probe for an ongoing attack, multiple successive control mode switches may be implemented on the process. To facilitate the detection of a wide range of attacks, successive switches from the nominal control mode to multiple attack-sensitive control modes may be considered. However, the closed-loop process under the switching-enabled detection strategy is a switched system. If not executed carefully, multiple control mode switches may destabilize the attack-free switched closed-loop system. To avoid potential instability of the attack-free closed-loop system, a minimum dwell time approach can be used whereby the process is forced to remain in each control mode for a minimum period of time before switching

to another control mode. The minimum dwell time required for closed-loop stability can be characterized using Lyapunov techniques for switched systems [34]. Utilizing the minimum dwell time approach to ensure the stability of the switched closed-loop system means that each switching instance is lower bounded by the minimum dwell time of process operation under the previous control mode. As a result, the switching instances under the randomized switching-enabled detection strategy are random subject to a constraint on stability.

While the dwell time for the nominal control mode ( $T_c^N$ ) may be chosen to maintain the stability of the attack-free switched closed-loop system, additional considerations may apply to the selection of a suitable dwell time for the process when operated under the attack-sensitive control mode ( $T_c^A$ ). One consideration is management of the tradeoff between attack detection and attack-free closed-loop performance. A longer dwell time in the attack-sensitive mode may increase the chances of attack detection, but it could also degrade closed-loop performance in the absence of attacks. Meeting process safety constraints could also be another consideration in the selection of the dwell time for the attack-sensitive mode. In the presence of an attack, the augmented state under the attack-sensitive control mode may grow unbounded with time. However, the augmented state may remain bounded within some safe set for a finite time of operation under the attack-sensitive mode. It is possible that the considerations for the selection of a suitable dwell time for operation in the attack-sensitive mode may be conflicting in some cases. For example, to ensure that process safety constraints are met, the dwell time for the attack-sensitive mode may be shorter than the minimum dwell time required for the stability of the attack-free switched closed-loop system. In such cases, the dwell time for the attack-sensitive mode should be chosen based on the process safety constraints, while the closed-loop stability of the attack-free switched closed-loop system can be guaranteed by implementing only a finite number of control mode switches over the infinite time interval.

Further design considerations for the switching-enabled active detection method are the switching instances ( $t_{s_i}$ ,  $i \in \{1, 2, 3, \dots\}$ ), which are randomly chosen time steps at which the control system switches from one control mode to the other. However, the stability considerations for the attack-free process may constrain control mode switching randomization. For switching from the nominal control mode to the attack-sensitive control mode, the switching instances are chosen as random integers such that the process is operated under the nominal mode for a time period that is at least equal to the dwell time for the nominal mode  $T_c^N$ . For switching from the attack-sensitive control mode to the nominal control mode, the switching instances are chosen such that the process is operated under the attack-sensitive control mode for the specified dwell time  $T_c^A$  (chosen to meet closed-loop stability, closed-loop performance, and process safety constraints). The first switching instance is a random number that is greater than or equal to the dwell time of the process operated under the nominal control mode, i.e.,  $t_{s_1} \geq T_c^N$ . All subsequent switching instances depend on the previous switching instance and account for the dwell time of the process under each mode. If no attack is detected, the even switching instances (for switching from the attack-sensitive control mode to the nominal control mode) may depend upon the prior switching instances and are chosen per the relation  $t_{s_{2n}} = t_{s_{2n-1}} + T_c^A$  for  $n \geq 1$ . Similarly, the switching instances for switching from the nominal mode to the attack-sensitive mode are selected as random integers that account for the dwell time under the nominal control mode per the relation  $t_{s_{2n+1}} \geq t_{s_{2n}} + T_c^N$  for  $n \geq 1$ .

Once an attack is detected, no further control mode switches are implemented in the process. From the results and discussion presented herein, monitoring the process using the reachable set-based detection scheme ensures that we can randomize control mode switches (randomly choose the switching instances  $t_{s_i}$ ) to enable attack detection while guaranteeing a zero false alarm rate. Randomization of the control mode switches may help preserve the confidentiality of the detection scheme from an attacker engaged in espionage and thereby prevent them from learning the switching schedule and designing a smart attack that could evade detection. Additionally, we do not make any assumptions on the

structure of the output matrix ( $C^x$ ). Therefore, the switching-enabled detection method present in this work may enable attack detection with a guaranteed zero false alarm rate even when  $C^x$  is non-invertible.

**Remark 1.** To select the attack-sensitive control parameters, the range of attacks to be detected may be chosen as attacks on the process operated under the nominal control mode which are either undetectable (attacks such that  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K^N, L^N)) \subseteq \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K^N, L^N))$  for all  $t \in \mathbb{Z}^+$ ) or potentially detectable (attacks such that  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K^N, L^N)) \not\subseteq \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K^N, L^N))$  for some  $t \in \mathbb{Z}^+$  and  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K^N, L^N)) \cap \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K^N, L^N)) \neq \emptyset$  for all  $t \in \mathbb{Z}^+$ ). More details on the selection of attack-sensitive control parameters may be found in [15].

**Remark 2.** The even switching instances may be selected randomly to vary the dwell time of the process operated under the attack-sensitive mode between the minimum dwell time that allows for the stability of the attack-free switched system ( $T_c^{A^{min}}$ ) and the specified dwell time chosen to meet the closed-loop performance, attack detection, stability, and process safety constraints ( $T_c^A$ ). Under this modified method, the even switching instances may be chosen as random integers that satisfy  $t_{s_{2n}} - t_{s_{2n-1}} \in [T_c^{A^{min}}, T_c^A]$  for all  $n \geq 1$ . However, implementing control mode switching with a time-varying dwell time for the attack-sensitive mode is subject to a rigorous characterization of the dwell times  $T_c^{A^{min}}$  and  $T_c^A$ , which is outside the scope of current work.

### 3.2. Algorithms for the Randomized Control Mode Switching-Enabled Cyberattack Detection Method

Algorithm 1 outlines the steps for the implementation of the randomized control mode switching-enabled cyberattack detection method. The proposed method considers that the reachable sets of the monitoring variable for the attack-free process are computed online and allows for multiple switches between the nominal and the attack-sensitive control modes to probe the process for attacks. The inputs to the algorithm are the set of initial states  $\mathcal{R}_0^{\xi}$ , the nominal control parameters ( $K^N, L^N$ ), the attack-sensitive control parameters ( $K^A, L^A$ ), the minimum dwell time for the attack-free process operated in the nominal control mode  $T_c^N$ , and the dwell time of the closed-loop process operated in the attack-sensitive control mode  $T_c^A$ . The detection time  $t_d$  is the output of the algorithm. The randomization of the control mode switching instances is dictated by the randomization flag  $f$ , which is a random variable that can take values of 0 or 1 and is assigned a new value at each time step (which is the controller sampling instance). If at a given time step,  $f = 0$ , then the current time step is not a switching instance, and if  $f = 1$ , then the current time step is a switching instance.

Algorithm 1 may be implemented with the reachable sets computed online using one of several methods that have been proposed in the literature (e.g., [35–37]). Online computation of the reachable sets may not scale well with the dimension of the state, and the computation may become intractable within the sampling instances. As a result, implementation of Algorithm 1 on a process under transient operation may not always be feasible.

Chemical processes are operated for extended periods at or near their steady states, where all possible values of the process states are bounded within the minimum invariant set. Considering this, we propose a modification to the algorithm for the implementation of control mode switching on a closed-loop process under steady-state operation using reachable sets computed offline. Based on standard results from the literature (see Theorem 1 in [33]), it can be shown that after a switch between the nominal and attack-sensitive control modes, the reachable sets of the attack-free process converge to an invariant neighborhood of the minimum invariant set of the process under the new mode in finite time. This means that, for the process at steady state, the values of the monitoring variable evolve within the terminal set of the operating mode prior to a control mode switch. Following each switching event, for the process under transient operation, its monitoring variable values

evolve within the attack-free reachable sets under the new mode. However, the transient operation lasts only for a finite number of time steps, until the reachable sets are contained entirely within (i.e., converge to) an invariant neighborhood of the minimum invariant set of the attack-free process operated under the new mode. After convergence, the process monitoring variable values are bounded within an invariant neighborhood of the terminal set of the attack-free process under the new mode. Based on these considerations, a hybrid approach may be used to monitor the switched closed-loop system. In this approach, the detection scheme switches with the control mode. The reachable set-based detection scheme is used only during the transient period after each switch. The detection scheme switches to a terminal set-based detection scheme after sufficient time has elapsed from the control mode switch such that the reachable sets of the augmented state for the attack-free process converge to an invariant neighborhood of the minimum invariant set of the process under the new mode. Using this hybrid monitoring approach, the computational load for computing the reachable sets during the implementation of Algorithm 1 may be reduced by terminating the online computation of the reachable sets after the process has attained steady-state operation (i.e., at the time step that the attack-free reachable sets converge to an invariant neighborhood of the minimum invariant set).

---

**Algorithm 1:** Algorithm for the randomized control mode switching-enabled attack detection with online reachable set computation.

---

**Inputs:**  $\mathcal{R}_0^\xi, (K^N, L^N), (K^A, L^A), T_c^N, T_c^A$   
**Initialization:**  $t_d = \infty, t_s = \infty, t = 0, (K_t, L_t) = (K^N, L^N)$   
**Outputs:**  $t_d$

```

1 do
2   Receive the measured variable  $y_t$  from the sensors over the sensor–controller
   communication link
3   Calculate the monitoring variable  $\eta_t$  and reachable set of the attack-free
   process  $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi)$ 
4   Randomization logic
5   if  $t_s = \infty$  then
6     Randomly generate switching flag :  $f \in \{0, 1\}$ 
7   else if  $t > t_s + T_c^A + T_c^N$  then
8     Set  $t_s = \infty$ 
9   Monitoring logic
10  Compute the detection scheme output  $\phi_t(\eta_t)$  per Equation (18)
11  if  $\phi_t(\eta_t) = 1$  then
12    Declare the detection of an attack. Set the detection time step to be the
    current time step  $t = t_d$ . Terminate the detection algorithm.
13  Switching logic
14  else if  $f = 1$  and  $t_s = \infty$  then
15    Set  $f = 0, t_s = t$  and  $(K_t, L_t) = (K^A, L^A)$ 
16  else if  $t = t_s + T_c^A$  then
17    Set  $(K_t, L_t) = (K^N, L^N)$ 
18  Set  $t \leftarrow t + 1, (K_{t+1}, L_{t+1}) = (K_t, L_t)$ 
19 while  $t_d = \infty$ ;

```

---

Algorithm 2 outlines the steps for implementing the randomized switching-enabled attack detection strategy for processes under steady-state operation that utilizes the hybrid monitoring approach. The inputs for the algorithm are the specified dwell time of the process operated under the nominal mode  $T_c^N$ , the dwell time of the process operated under the attack-sensitive mode  $T_c^A$ , the nominal  $(K^N, L^N)$  and the attack-sensitive  $(K^A, L^A)$  control parameters, the minimum invariant sets for the attack-free process operated under

the nominal mode  $\mathcal{R}_\infty^\xi(K^N, L^N)'$  and under the attack-sensitive mode  $\mathcal{R}_\infty^\xi(K^A, L^A)'$ , the number of time steps that the reachable sets take to converge to the minimum invariant set for the process operated under the attack-sensitive mode (starting from the minimum invariant set for the process operated under the nominal mode  $t_r$ ), the number of time steps the reachable sets take to converge to the nominal minimum invariant set (starting from the attack-sensitive minimum invariant set  $t_r^*$ ), and the reachable sets for a switch from the nominal control mode to the attack-sensitive control mode  $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K^A, L^A))$  for  $t \in (0, t_r]$  and for a switch from the attack-sensitive control mode to the nominal control mode  $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K^N, L^N))$  for  $t \in (0, t_r^*]$ . The algorithm implementation is terminated when an attack is detected.

---

**Algorithm 2:** Algorithm for the randomized control mode switching-enabled attack detection for processes with offline computation of reachable sets.

---

**Inputs:**  $\mathcal{R}_\infty^\eta(\mathcal{R}_\infty^\xi(K^N, L^N)'), \mathcal{R}_\infty^\eta(\mathcal{R}_\infty^\xi(K^A, L^A)'), T_c^A, T_c^N,$   
 $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K^A, L^A))$  for  $t \in (0, t_r], \mathcal{R}_t^\eta(\mathcal{R}_t^\xi(K^N, L^N))$  for  $t \in (0, t_r^*]$   
**Initialization:**  $t = 0, t_d = \infty, t_s = \infty, \mathcal{R}_0^\xi = \mathcal{R}_\infty^\xi(K^N, L^N)', (K_t, L_t) = (K^N, L^N)$   
**Outputs:**  $t_d$

```

1 do
2   Receive the measured variable  $y_t$  from the sensors over the sensor-controller
   communication link
3   Calculate the monitoring variable  $\eta_t$ 
4   Randomization logic
5   if  $t_s = \infty$  then
6     Randomly generate switching flag :  $f \in \{0, 1\}$ 
7   else if  $t > t_s + T_c^A + T_c^N$  then
8     Set  $t_s = \infty$ 
9   Reachable sets for monitoring
10  if  $t_s = \infty$  or  $t > t_s + t_r^*$  then
11     $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi) = \mathcal{R}_\infty^\eta(\mathcal{R}_\infty^\xi(K^N, L^N)')$ 
12  else if  $t \leq t_s + t_r$  then
13     $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi) = \mathcal{R}_p^\eta(\mathcal{R}_t^\xi(K^A, L^A))$ 
14     $p = t + 1 - (t_s + t_r)$ 
15  else if  $t \in [t_s + T_c^A, t_s + T_c^A + t_r^*)$  then
16     $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi) = \mathcal{R}_p^\eta(\mathcal{R}_t^\xi(K^N, L^N))$ 
17     $p = t + 1 - (t_s + T_c^A + t_r^*)$ 
18  else if  $t \in [t_s + t_r, t_s + T_c^A)$  then
19     $\mathcal{R}_t^\eta(\mathcal{R}_t^\xi) = \mathcal{R}_\infty^\eta(\mathcal{R}_\infty^\xi(K^A, L^A)')$ 
20  Monitoring logic
21  Compute the detection scheme output  $\phi_t(\eta_t)$  per Equation (18)
22  if  $\phi_t(\eta_t) = 1$  then
23    Declare the detection of an attack. Set the detection time step to be the
    current time step  $t = t_d$ . Terminate the detection algorithm.
24  Switching logic
25  else if  $f = 1$  and  $t_s = \infty$  then
26    Set  $f = 0, t_s = t$  and  $(K_t, L_t) = (K^A, L^A)$ 
27  else if  $t = t_s + T_c^A$  then
28    Set  $(K_t, L_t) = (K^N, L^N)$ 
29  Set  $t \leftarrow t + 1, (K_{t+1}, L_{t+1}) = (K_t, L_t)$ 
30 while  $t_d = \infty;$ 

```

---



**Remark 3.** Per the steps in Algorithms 1 and 2, the odd switching instances, which dictate when a switch from the nominal control mode to attack-sensitive control mode occurs, may be chosen as an arbitrarily large positive integer. To avoid waiting for an inordinately long time period before a control mode switch, the time step for odd switching instances may be chosen as an integer over a finite time interval  $[t_s^{\min}, t_s^{\max}]$ , where  $t_s^{\min} \geq T_c^N$  is the lower bound and  $t_s^{\max} \geq t_s^{\min}$  is the operator-specified upper bound of the interval.

**Remark 4.** When implementing Algorithms 1 and 2, the randomization flag ( $f$ ) may be drawn from a Bernoulli distribution with  $p \in [0, 1]$  and  $q = 1 - p$ . Where  $p$  is the probability that a random variable drawn from the distribution takes a value of 1, and  $q$  is the probability that a random variable from the distribution takes a value of 0. If  $p = 0.5$ , the process is switched as frequently as it is not switched. To increase the likelihood of confounding an attacker,  $p$  may be chosen as a number that is not equal to 0.5.

#### 4. Application of Randomized Control Mode Switching Algorithms to Illustrative Processes

In this section, we demonstrate the application of the proposed switching-enabled attack detection strategy with randomized control mode switching using two illustrative examples. In the first example, we demonstrate the application of Algorithm 1 on a dynamic process and use attack-free reachable sets computed online for monitoring the process. In the second example, we demonstrate the application of Algorithm 2 to a chemical process example under steady-state operation and use attack-free reachable sets computed offline to monitor the process. Using simulations of the chemical process example, we also demonstrate the application of the randomized control mode switching approach for the detection of a “smart” attack that is designed to evade detection under a scheduled control mode switch. The MPT 3.0 toolbox [38] is used to compute all polytopes [38] and the CORA toolbox [39] is used to compute all zonotopes.

##### 4.1. Illustrative Scalar Process

We consider an example process that has one state and one manipulated input:

$$\begin{aligned}x_{t+1} &= x_t + u_t + w_t \\y_t &= \Lambda^y(x_t + v_t) \\u_t &= \Lambda^u(-K\hat{x}_t)\end{aligned}\tag{20}$$

where  $x_t \in \mathbb{R}$ ,  $y_t \in \mathbb{R}$ , and  $u_t \in \mathbb{R}$  for all  $t \in \mathbb{Z}^+$  are the process state, the measured output, and the manipulated input, respectively. The scalars,  $v_t \in \mathcal{V} := \{v \in \mathbb{R} \mid |v| \leq 1\}$  and  $w_t \in \mathcal{W} := \{w \in \mathbb{R} \mid |w| \leq 1\}$ , model the bounded measurement noise and the process disturbance, respectively. The process may be under a multiplicative attack that modifies the operational data over all PCS communication links, which are represented by  $\Lambda^y \neq 1$  and  $\Lambda^u \neq 1$ . To analyze the closed-loop process, we define an augmented state that is a concatenation of the process state and the estimation error as follows:  $\xi := [x \ e]^T$ . With this definition, the evolution of the augmented state of the process is described by Equation (9) with  $A^x = 1$ ,  $B^u = 1$ ,  $C^x = 1$ . For process monitoring, we consider a monitoring variable that is a concatenation of the measured output and the residual vectors  $\eta := [y \ r]^T$ . The monitoring variable  $\eta$  may be expressed in the form of in Equation (12) with  $H^y = [1 \ 1]^T$  and  $H^g = [0 \ -1]^T$ .

We choose the nominal controller gain  $K^N$  to minimize the quadratic cost  $J = \mathbb{E}[\sum_{i=0}^{\infty} (x_i^T Q x_i + u_i^T R u_i)]$ , with  $Q = 5$ ,  $R = 1$ , and the nominal observer gain as the steady-state Kalman filter gain with covariance matrices  $Q_L = 0.0011$  and  $R_L = 0.0011$ . Likewise, we choose the attack-sensitive parameters such that the attacked closed-loop process operated under the attack-sensitive mode is unstable in the sense that  $\rho(A^{\xi_a}(K^A, L^A)) > 1$  over the attack range  $\Lambda^y = 0.86$  and  $\Lambda^u \in [1.1, 2]$ , and  $\Lambda^u = 1.1$  and  $\Lambda^y \in [0.1, 0.99]$ . The matrix pair  $(A^{\xi_a}(K^A, L^A), C^{\eta_a})$  is observable over the attack range considered for selecting

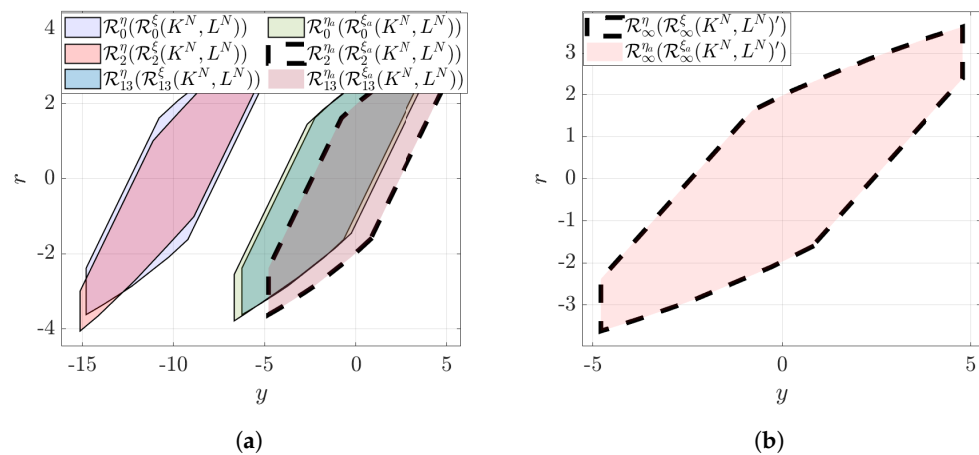


the attack-sensitive parameters. The values of the nominal and attack-sensitive control parameters are  $(K^N, L^N) = (0.8541, 0.618)$  and  $(K^A, L^A) = (1.57, 1.28)$ . For the attack-free process, we compute invariant outer  $\epsilon$ -approximations of the minimum invariant set for the process under the nominal and the attack-sensitive modes with an error bound of  $\epsilon = 5 \times 10^{-5}$  using the method described in [32]. For brevity, in the remainder of this paper, we refer to the invariant outer  $\epsilon$ -approximation of the minimum invariant set as the minimum invariant set. For the attack-free process with the control system under the nominal mode, its reachable sets take 18 time steps to converge from the set of initial states to the minimum invariant set. Similarly, for the attack-free process operated exclusively under the attack-sensitive mode, the reachable sets of the augmented state from the set of initial states converge to the minimum invariant set in five time steps.

We consider the process under a transient operation when its state evolves from a set of initial states, that is, the polytope obtained by shifting all the vertices of the minimum invariant set of the attack-free process operated under the nominal mode by  $\xi' = [-10 \ 0]^T$ . To quantify the performance of the controller for the attack-free process operated under the attack-sensitive mode and the nominal mode, we perform two sets of simulations (each set consisting of 1000 simulations) of the attack-free process. In the first set, we consider the exclusive operation of the attack-free process under the nominal mode. In the second set, we consider the exclusive operation of the attack-free process under the attack-sensitive mode. Within a simulation set, at each time step of each simulation, we vary the values of the process disturbance and the measurement noise, which we model as random variables drawn from two separate normal distributions with  $\mathcal{N}(0, 0.0333)$ . However, across simulation sets, we use the same values of process disturbance and measurement noise. Within each simulation, we consider the evolution of the process for 1000 time steps, and we initialize the state at  $\xi_0 = [-10 \ 0]^T$ . We compare the quadratic cost ( $J = \mathbb{E}[\sum_{i=0}^{1000} (x_i^T Q x_i + u_i^T R u_i)]$ ) across the simulation sets. Over the simulations of the process operated under the nominal mode, the average quadratic cost was found to be 648.66 with a standard deviation of 37.59. Similarly, over simulations of the process operated under the attack-sensitive mode, the average quadratic cost was found to be 1189.68 with a standard deviation of 156.77. Comparing the performance of the controller between the two modes, it can be concluded that the controller performance is worse under the attack-sensitive mode. This result demonstrates that to manage the tradeoff between attack detection and the attack-free performance degradation resulting from control mode switching, intermittent switching from the nominal control mode to the attack-sensitive control mode may be preferred to operation under the attack-sensitive control mode exclusively.

We now apply the switching-enabled detection method (Algorithm 1) over simulations of the process during the transient operation for the detection of an attack with  $\Lambda^y = 0.86$  and  $\Lambda^u = 1.1$ . First, we analyze the detectability of the attack when the process is operated under the nominal control mode by comparing the reachable sets associated with the attack-free process to those associated with the attacked process. Figure 1a illustrates the reachable sets for the process operated under the nominal mode over a few time steps during transient operation. As illustrated, at time steps  $t = 0, 2, 13$ , the reachable sets for the attacked and the attack-free process always intersect; however, the attacked reachable sets are not contained within the attack-free reachable sets, meaning that the attack is potentially detectable. The transient operation of the attacked process lasts over 12 time steps over the time interval  $t \in [0, 13)$ . Over all time steps during transient operation of the process, the attack is found to be potentially detectable because the reachable sets satisfy  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K^N, L^N)) \cap \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K^N, L^N)) \neq \emptyset$  and  $\mathcal{R}_t^{\eta_a}(\mathcal{R}_t^{\xi_a}(K^N, L^N)) \not\subseteq \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K^N, L^N))$ . At time step  $t = 13$ , the reachable sets of the attacked process converge to the terminal set of the attacked process. Therefore, to analyze attack detectability over the time steps  $t \in (13, 18]$ , the terminal set of the process under attack is compared with the reachable sets of the attack-free process, and the attack is found to be undetectable. Figure 1b illustrates the terminal set of the attack-free process and the terminal set of the attacked process,

showing that the attack on the process under steady-state operation is undetectable due to the fact that the terminal set of the attacked process is contained entirely within the terminal set of the attack-free process, i.e.,  $\mathcal{R}_t^{\eta^a}(\mathcal{R}_t^{\xi_a}(K^N, L^N)) \subset \mathcal{R}_t^{\eta}(\mathcal{R}_t^{\xi}(K^N, L^N))$ . To verify attack detectability analysis, we consider 1000 simulations of the attacked process operated exclusively under the nominal mode and monitored by the reachable set-based detection scheme in Equation (15). We initialize each simulation at  $\xi_0 = [-10 \ 0]^T$  and consider the evolution of the process states over 1000 time steps. Over all simulations, we vary the values of the process disturbances and measurement noise similar to the previous simulation sets. Over all simulations considering the process under an attack, the attack is not detected.



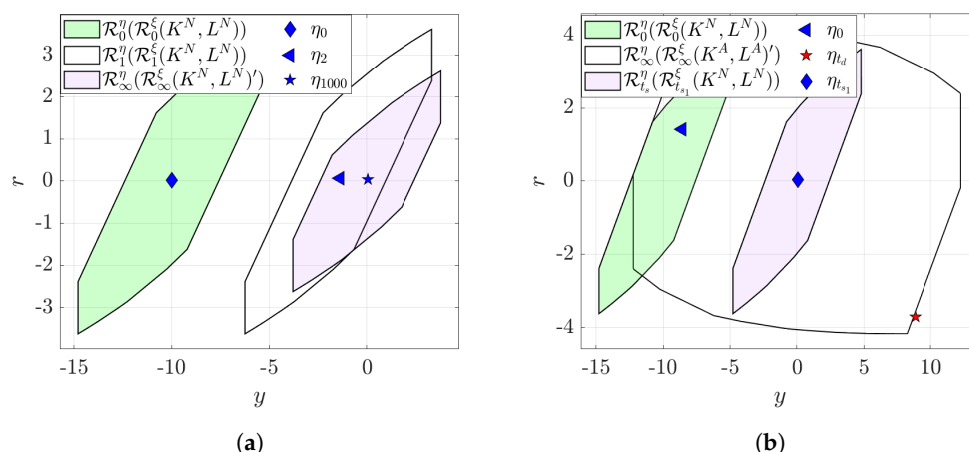
**Figure 1.** The reachable sets of the attacked and the attack-free process operated under the nominal mode (a) during transient operation when the attack is potentially detectable, and (b) at the steady state when the attack is undetectable with  $\mathcal{R}_{\infty}^{\eta^a}(\mathcal{R}_{\infty}^{\xi_a}(K^N, L^N)) \subset \mathcal{R}_{\infty}^{\eta}(\mathcal{R}_{\infty}^{\xi}(K^N, L^N))$ .

We apply the switching-enabled active detection method over simulations of the transient closed-loop process, and for monitoring, we use the reachable set-based detection scheme in Equation (15). We design two sets of simulations (each consisting of 1000 simulations of the process) similar to the simulations considered earlier when comparing the quadratic cost under the nominal and attack-sensitive modes. In the first simulation set, we consider the attack-free process, while in the second simulation set, we consider the attacked process with the attack beginning at time step  $t = 0$ . In this section, detection of an attack on the process during transient operation is considered. Over each simulation, the switching-enabled attack detection method implements a single control mode switch at a randomly chosen time instance in the interval  $[0, 17]$  when the attack-free process states under the nominal control mode are not within the minimum invariant set for the process operated under the nominal mode. Since the simulations consider a single switch between the two modes, we do not specify a minimum dwell time for each mode. Similarly, over simulations that consider the process operated under the attack-sensitive control mode, we use a dwell time of  $T_c^A = 150$ . Across simulations of the attack-free and the attacked processes, we consider the same switching instance. To implement the detection scheme, we compute the attack-free reachable sets online at each time step by using Equation (14) with  $(K, L)$  selected based on the process operation mode (under the nominal mode  $(K, L) = (K^N, L^N)$ , while under the attack-sensitive mode  $(K, L) = (K^A, L^A)$ ). However, to reduce the computational load, the online computation of the reachable sets is terminated at the time step when the attack-free augmented state is expected to be contained within the minimum invariant set for the process operated under the mode of operation considered. After termination of online computation of the reachable sets, the terminal set of the monitoring variable for the attack-free process under the mode of operation considered is used to monitor the process.

Over the 1000 simulations that consider the attack-free process, a switch from the nominal to the attack-sensitive control mode is implemented at time steps chosen randomly over the time interval  $[0, 17]$ . No attack detection occurs over all simulations after the first switch from the nominal control mode to the attack-sensitive control mode, and the control mode switches back to the attack-sensitive mode after 150 time steps from the first switch. No false alarms were observed over all simulations. Under the second simulation set, detection of the attack occurs over all simulations within a minimum of 5 time steps and a maximum of 135 time steps from the control mode switch. Over all simulations, the control system switches back to the nominal control mode, after which the attack is not detected. Figure 2a,b show the values of the monitoring variable over a few time steps of one simulation considering the attack-free process, and one simulation considering the attacked process, respectively. Over both simulations, the control parameters switch from the nominal to attack-sensitive values at time step  $t_{s_1} = 14$ .

Figure 2a shows the values of the monitoring variable values observed when the switch is implemented over a simulation considering the attack-free process. At time step  $t = 0$ , no false alarm is generated as the value of the monitoring variable at that time step represented by the blue diamond marker is contained within the attack-free reachable set at that time step, which is the set shown in green. Similarly, no false alarm is observed at time step  $t = 1$  because the monitoring variable value represented as the blue diamond is contained within the attack-free reachable set at that time shown as the set in white. While omitted for clarity, no false alarms are observed until the control mode switch at  $t_{s_1} = 14$  because the monitoring variable values at each time step evolve within the corresponding attack-free reachable sets. After the switch is implemented, the process is operated under the attack-sensitive mode for 150 time steps, during which no false alarms are observed. At time step  $t_{s_1} + T_c^A = 164$ , an attack is not detected, causing the control system to switch back to the nominal mode. No false alarms are observed even after this switch until the end simulation, at time step  $t = 1000$ , when the monitoring variable value represented by the blue star marker is contained within the attack-free terminal set shown as the set in purple.

Figure 2b shows the values of the monitoring variable over a few time steps of a simulation considering the process under the attack. At time step  $t = 0$ , the attack on the process under the nominal mode is not detected because the monitoring variable value represented by the blue triangle marker is contained within the attack-free reachable set at that time step, shown as the green set. While omitted for clarity, the attack is not detected during the process operation under the nominal mode because the monitoring variable values over the time interval  $t \in [0, 14)$  are contained within the corresponding attack-free reachable sets. After a switch from the nominal control mode to the attack-sensitive control mode, the attack is detected the time step  $t_d = 124$  because the monitoring variable value represented by the red star marker leaves the attack-free reachable set (which is the terminal set of the process operated under the attack-sensitive control mode) at that time step shown as the white set. In this case, attack detection occurs after the reachable sets of the attack-free process converge to its terminal set under the attack-sensitive mode. After the attack is detected, the control systems switches back to the nominal control mode and no further alarms are observed because the monitoring variable values at each time step are contained within the corresponding attack-free reachable sets until the end of simulation when the attack-free reachable sets have converged to the terminal set of the attack-free process under the nominal control mode. The monitoring variable value at time step  $t = 1000$  (end of the simulation) is represented by the blue diamond marker, and as shown, no alarm is observed at this time step because the monitoring variable is contained within the attack-free terminal set shown as the set in purple. These results demonstrate that the switching-enabled attack detection method utilizing the reachable set-based detection scheme enables attack detection on a dynamic process while guaranteeing a zero false alarm rate due to a control mode switch implemented at a randomly chosen time step.



**Figure 2.** Evolution of values of the monitoring variable, with respect to the attack-free reachable sets, over a few time steps of a simulation of the scalar process with a control mode switch implemented at time step  $t_{s1} = 14$  for the case when (a) no attack takes place (demonstrating zero false alarms), and (b) an attack takes place at  $t = 0$  and is detected at time step  $t_d = 124$ .

#### 4.2. A Continuous Stirred Tank Reactor

We consider an example process that consists of a continuous stirred tank reactor (CSTR) with a second-order exothermic reaction of the form  $A \rightarrow B$ . The process dynamics are modeled by the following system of ordinary differential equations:

$$\begin{aligned} \frac{dC_A}{dt} &= \frac{F}{V}(C_{A0} + \Delta C_{A0} - C_A) - k_0 e^{\frac{-E}{RT}} C_A^2 \\ \frac{dT}{dt} &= \frac{F}{V}(T_0 + \Delta T_0 - T) - \frac{\Delta H k_0}{\rho C_p} e^{\frac{-E}{RT}} C_A^2 + \frac{Q}{\rho C_p V} \end{aligned} \quad (21)$$

where  $C_{A0}$  and  $T_0$  are the inlet reactant concentration and feed temperature, respectively, and  $C_A$  and  $T$  are the reactant concentration and reactor temperature, respectively. The rate of heat transfer to or from the reactor  $Q$  is chosen as the manipulated input. The process is subject to bounded disturbances modeled as variations in the inlet reactant concentration  $\Delta C_{A0}$  and variations in the feed temperature  $\Delta T_0$ . The bounded process disturbances are within the limits  $|\Delta C_{A0}| \leq 0.015 \text{ kmol m}^{-3}$  and  $|\Delta T_0| \leq 4.5 \text{ K}$ . The measured variable available to the controller is the reactor temperature  $T$ . The bounded noise in the measurements from the sensor is within limits such that  $|v| \leq 4.5 \text{ K}$ . Table 1 provides a list of the definitions and values of the process parameters. Because the measurement of all possible states of the reactor are not available to the controller, the illustrative example presented in this section considers a case where the output matrix  $C$  is non-square and non-invertible.

To obtain a model similar to Equation (1), the continuous-time nonlinear process model in Equation (21) is discretized using a sampling interval of  $\Delta = 1 \times 10^{-2} \text{ h}$ . The system matrices for the linearized CSTR process are:

$$A^x = \begin{bmatrix} 0.7364 & -0.0041 \\ 10.6953 & 1.156 \end{bmatrix}, B^u = \begin{bmatrix} -0.0009 \times 10^{-4} \\ 0.4674 \end{bmatrix}, B^w = \begin{bmatrix} 0.0433 & -0.001 \\ 0.2724 & 0.054 \end{bmatrix}$$

The nominal observer gain is chosen as the steady-state Kalman filter gain with covariance matrices,  $Q_K = \begin{bmatrix} 1.4062 \times 10^{-5} & 0 \\ 0 & 1.2656 \end{bmatrix}$ ,  $R_K = 1.2656$  and the nominal controller gain is chosen to minimize the quadratic cost  $J = \mathbb{E}[\sum_{i=0}^{\infty} (x_i^T Q x_i + u_i^T R u_i)]$  with  $Q = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$  and  $R = 10$ . The attack-sensitive control parameters are chosen so that the attacked closed-loop process under the attack-sensitive mode is unstable under a range of attacks by checking if  $\rho(A^{\xi_a}(K^A, L^A)) > 1$  over the attack range  $\Lambda^y = 1.1$ ,  $\Lambda^u \in [1.3, 2]$ , and

$\Lambda^u = 1.3, \Lambda^y \in [1.1, 2]$ . Over the attack range considered, the matrix pair  $(A^{\zeta_a}(K^A, L^A), C^{\eta_a})$  is observable, indicating that all attacks in the range considered are potentially detectable.

**Table 1.** CSTR process parameters.

Flow rate ( $F$ )	$5.0 \text{ m}^3 \text{ h}^{-1}$
Reactor volume ( $V$ )	$1.0 \text{ m}^3$
Feed concentration of $A$ ( $C_{A0}$ )	$4.0 \text{ kmol m}^{-3}$
Activation energy ( $E$ )	$5.0 \times 10^4 \text{ kJ kmol}^{-1}$
Pre-exponential factor ( $k_0$ )	$8.46 \times 10^6 \text{ m}^3 \text{ h}^{-1} \text{ kmol}^{-1}$
Gas constant ( $R$ )	$8.314 \text{ kJ kmol}^{-1} \text{ K}$
Feed temperature ( $T_0$ )	$300 \text{ K}$
Density of reaction mass in the reactor ( $\rho$ )	$1000 \text{ kg m}^{-3}$
Heat of reaction ( $\Delta H$ )	$-1.15 \times 10^4 \text{ kJ kmol}^{-1}$
Heat capacity ( $C_p$ )	$0.231 \text{ kJ kg K}^{-1}$
Heat rate added/removed from the CSTR at the steady state ( $Q_s$ )	$0 \text{ kJ h}^{-1}$
Steady-state reactant concentration ( $C_{As}$ )	$1.22 \text{ kmol m}^{-3}$
Steady-state temperature ( $T_s$ )	$438.2 \text{ K}$

In this section, we apply the switching-enabled detection method on the CSTR process under steady-state operation, when all values of the process states are bounded within the minimum invariant set of the process. The disturbance set is modeled as a zonotope with the origin as the center. To compute invariant outer  $\epsilon$ -approximations of the minimum invariant set (henceforth referred to as the minimum invariant set for brevity) of the attack-free process under the nominal and the attack-sensitive modes, we use the method in [40] and use an error bound of  $\epsilon = 5 \times 10^{-5}$ . For a switch performed on the attack-free process from the nominal control mode to the attack-sensitive control mode, we compute the reachable sets with the set of initial states  $\mathcal{R}_0^{\zeta} = \mathcal{R}_{\infty}^{\zeta}(K^N, L^N)'$  for  $t_r = 11$  time steps when the sets are contained entirely within the minimum invariant set of the attack-free process under the attack-sensitive mode. Similarly, for a switch from the attack-sensitive control mode to the nominal control mode, we compute the reachable sets for the attack-free process with the set of initial states  $\mathcal{R}_0^{\tilde{\zeta}} = \mathcal{R}_{\infty}^{\tilde{\zeta}}(K^A, L^A)'$  until they converge to the minimum invariant set under the nominal control mode after  $t_r^* = 150$  time steps. We use the monitoring variable vector  $\eta = [y \ r]^T$  and compute its terminal set under each control mode and the reachable sets after each control mode switch until they are contained entirely within the terminal set under the new mode.

Initially, we perform two simulation sets to compare the quadratic cost ( $J = \mathbb{E}[\sum_{i=0}^{1500} (x_i^T Q x_i + u_i^T R u_i)]$ ) for operating the process exclusively under the nominal mode with the cost of operation of the system exclusively under the attack-sensitive mode. The first set of simulations consider the attack-free process operated exclusively under the nominal mode, while the second set of simulations considers the exclusive operation of the attack-free process under the attack-sensitive mode. Under each set, 1000 simulations considering the process operating under attack-free conditions are conducted. Each simulation considers the evolution of the process over 1500 time steps spanning 15 h in real time. We model the process disturbances on the feed concentration as random variables drawn from a distribution with  $\mathcal{N}(0, 0.005)$ , and we model the process disturbance and measurement noise on the feed temperature and the temperature of the reactor (the measured output) as random variables picked from two distinct normal distributions with  $\mathcal{N}(0, 1.5)$ . Over each time step of each simulation, we vary the values of the random variables representing the process distribution and the measurement noise. However, we use the same values of random numbers across simulation sets. The disturbances are clipped at the absolute value of their bound to ensure that there are no false alarms (e.g., if the absolute value random number representing  $\Delta C_{A0}$  exceeds 0.015, it is set to 0.015). Each simulation is initialized at the origin, which is contained within the minimum invariant sets of the process under



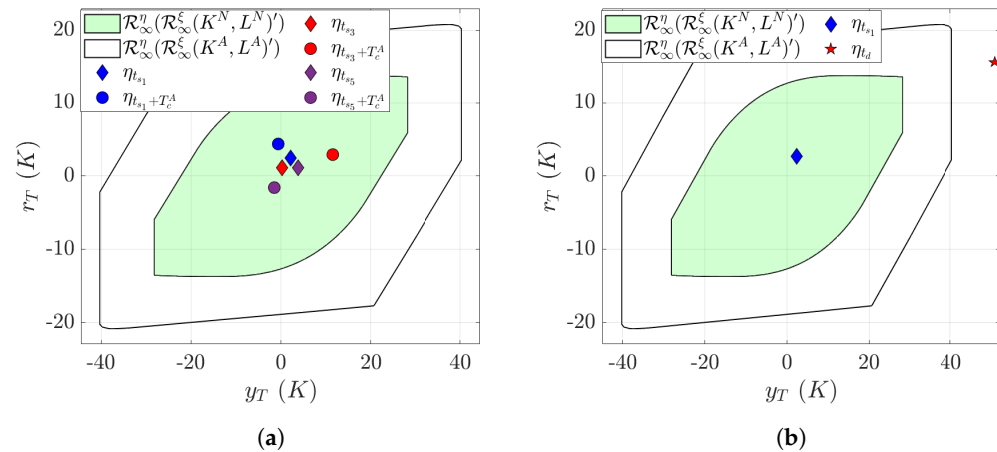
the attack-sensitive and the nominal modes. Over all simulations, considering the process operated under the nominal mode, the quadratic cost has a mean of  $1.016 \times 10^4$  and a standard deviation of  $2.498 \times 10^3$ . However, for the process operated exclusively under the attack-sensitive mode, the quadratic cost has a mean of  $1.1985 \times 10^{15}$  and a standard deviation of  $6.2674 \times 10^{13}$ . Therefore, the performance of the controller under the attack-sensitive mode is higher than under the nominal mode, indicating that switching may be preferable to extended process operation under the attack-sensitive control mode.

Next, we perform two simulation sets of the process monitored by the reachable set-based detection scheme in Equation (15). In the first set, we implement Algorithm 2 over simulations considering the process without an attack to show that there are no alarms generated from consecutive control mode switches implemented at randomly chosen time steps. To this end, we consider the detection of an attack with  $\Lambda^y = 1.1$  and  $\Lambda^u = 1.3$ . The attack beginning at time step  $t = 0$  on the process operated only under the nominal mode is not detected over 1000 simulations, even though it is potentially detectable. Therefore, in the second simulation set, we implement Algorithm 2 over 1000 simulations of the attacked process (attack begins at  $t = 0$ ). Over simulations considering the process operated under the attack-sensitive mode, we restrict the dwell time under attack-sensitive mode to  $T_c^A = 150$ . Over each simulation, the switching instances are chosen randomly such that a maximum of three switches from the nominal to the attack-sensitive control mode and back from the attack-sensitive to the nominal control mode are possible over each simulation. Specifically, the first switching instance from the nominal to the attack-sensitive control mode ( $t_{s_1}$ ) is selected as a random integer in the interval  $[0, 1200]$ . The second switching instance from the nominal control mode ( $t_{s_3}$ ) is based on the first switching instance by selecting a random integer over the interval  $[t_{s_1} + T_c^A + 150, 1200]$ . If no attack is detected, then we allow for a third instance of a switch from the nominal control mode to the attack-sensitive control mode, with the switching instance  $t_{s_5}$  chosen as a random integer over the interval  $[t_{s_3} + T_c^A + 150, 1200]$ . We do not specify a minimum dwell time for operation under the nominal mode since the simulations consider a finite number of switches between the different control modes.

Over all simulations considering the attack-free process, no alarms are observed, and a minimum of two and a maximum of three control mode switches are implemented. Over all simulations, the process is under the nominal mode at the end of the simulation. Figure 3a illustrates the values of the monitoring variable observed over a few time steps of a simulation of the attack-free process over which there are three consecutive switches from the nominal to the attack-sensitive control mode implemented at time steps  $t_{s_1} = 221$ ,  $t_{s_3} = 767$ , and  $t_{s_5} = 1313$ . Over all simulations of the attack-free process, no false alarms are observed. While omitted for clarity in Figure 3a, the monitoring variable values before the time step  $t_{s_1}$  evolve within the terminal set of the attack-free process under the nominal control mode shown as the green set and no false alarms are observed until the first switching instance  $t_{s_1} = 221$  when the control system switches to the attack-sensitive control mode. No false alarms are observed during process operation under the attack-sensitive control mode. As a result, after the dwell time under the attack-sensitive mode elapses, the control system switches back to the nominal mode at time step  $t_{s_1} + T_c^A$ . As shown, no false alarm is observed at  $t_{s_1} + T_c^A$  because the monitoring variable value (indicated by the blue circle marker) is contained within the terminal set of the attack-free process under the attack-sensitive control mode (indicated by the white set). While omitted here for brevity, no false alarms are observed during process operation under the nominal control mode, and a second switch to the attack-sensitive control mode occurs at time step  $t_{s_2} = 767$ , when the monitoring variable value (indicated by the red diamond marker) is contained within the attack-free terminal set. Even after the second control mode switch, no false alarms are observed and the control system switches back to the nominal control mode after the dwell time under the attack-sensitive mode elapses at time step  $t_{s_2} + T_c^A$ , when the monitoring variable value (indicated by the red circle marker) is contained within the terminal set of the process under attack-sensitive mode. No false alarms are observed



even after a third switch from the nominal to the attack-sensitive control mode occurs at  $t_{s_3}$  (monitoring variable value shown by purple diamond marker) followed by a switch back from the attack-sensitive to the nominal control mode at time step  $t_{s_3} + T_c^A$  (monitoring variable value shown by purple circle marker). The results demonstrate that the proposed control mode switching strategy guarantees a zero false alarm rate when implemented on the attack-free process.



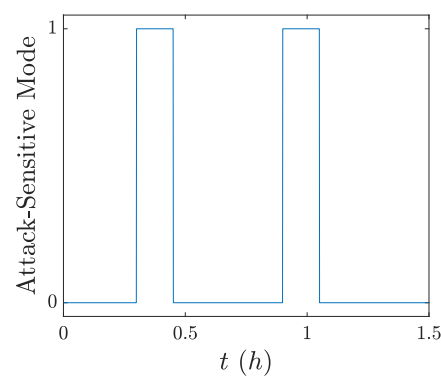
**Figure 3.** Values of the monitoring variable observed over a few time steps of a simulation considering (a) the attack-free process with consecutive switches from the nominal control mode to the attack-sensitive control mode implemented at  $t_{s_1} = 221$ ,  $t_{s_3} = 767$ , and  $t_{s_5} = 1313$  showing no false alarms due to switching and, (b) the process under an attack, with switching from the nominal to the attack-sensitive control mode at time step  $t_{s_1} = 221$  leading to attack detection at time step  $t_d = 235$ .

Figure 3b illustrates the values of the monitoring variable observed over some time steps of one simulation of the process under the attack, over which the first switching instance is  $t_{s_1} = 221$ . Until this control mode switch, the attack is not detected because the monitoring variable values evolve within the terminal set of the attack-free process under the nominal control mode. At the switching instance  $t_{s_1}$ , the monitoring variable value shown by the blue marker is contained within the terminal set of attack-free process under the nominal mode (shown as the green set), meaning that no attack is detected. After the control mode switch, no alarms are observed until the attack is detected at time step  $t_d = 235$  when the reachable sets of the attack-free process are contained entirely within the terminal set of the attack-free process operated under the attack-sensitive control mode. As shown in Figure 3b, at the detection time step  $t_d$ , the monitoring variable value of the process shown by the red diamond marker is not contained within the terminal set of the attack-free process operated under the attack-sensitive control mode. Over all simulations of the process under attack, detection of the attack occurs after the first control mode switch within a minimum of 4 time steps and a maximum of 69 time steps from the switching instance. The results demonstrate that the reachable set-based detection scheme in Equation (15) guarantees attack detection with a zero false alarm rate, for a randomly chosen switching instance, even when the output matrix  $C^x$  is non-invertible.

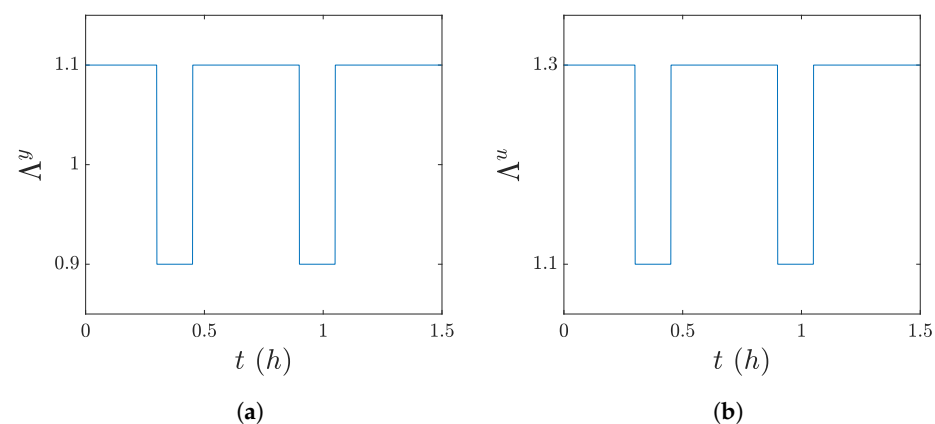
#### 4.2.1. Comparison of Randomized and Scheduled Control Mode Switching

In this section, we demonstrate the application of the randomized control mode switching to enable the detection of a “smart” attack that is designed to evade detection under a scheduled control mode switching-enabled attack detection method. We consider a simulation of the CSTR process, over which control mode switching is implemented per a fixed schedule as shown in Figure 4. As illustrated, the switching schedule allows for two control mode switches between the nominal and the attack-sensitive control modes. For switching from the nominal to the attack-sensitive control mode, the first switching instance is  $t_{s_1} = 300$  and the second switching instance is  $t_{s_3} = 900$ . If no attack is detected

on the process operated under the attack-sensitive control mode until the dwell time of  $T_c^A = 150$  time steps elapses, the controller switches back to the nominal mode. We assume that an attacker who is aware of the switching schedule designs a smart attack that switches the attack matrices in sync with the control mode switching. Specifically, the attacker uses an attack with  $\Lambda^y = 1.1$  and  $\Lambda^u = 1.3$  when the process is expected to operate under the nominal control mode. While potentially detectable, the attack on the process operated under the nominal control mode is not detected. However, when the process is expected to operate under the attack-sensitive control mode, the attacker switches to using an attack with  $\Lambda^y = 0.9$  and  $\Lambda^u = 1.1$  under which the attacked process is stable. Similar to the attack on the process operated under the nominal mode, the attack on the process operated under the attack-sensitive mode is potentially detectable; however, the attack is such that it is not detected. The attack schedule is as shown in Figure 5a,b.

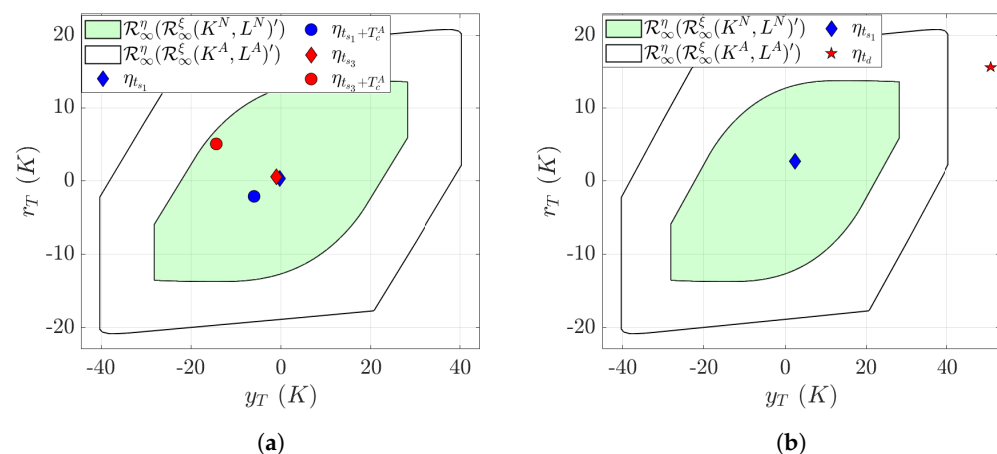


**Figure 4.** Scheduled control mode switching for the CSTR process. A value of 0 on the Y-axis indicates that the process is operated under the nominal control mode, and a value of 1 on the Y-axis indicates that the process is operated under the attack-sensitive control mode.



**Figure 5.** The smart attack schedule, which is designed to evade detection under the scheduled control mode switching-based detection scheme. The figures illustrate the attack schedule for (a) the sensor–controller attack and (b) the controller–actuator attack.

We consider one simulation that considers the closed-loop process under the smart attack, with the control mode switch implemented per the schedule shown in Figure 4. The process is initialized under the nominal control mode, with its initial state chosen as the origin. The reachable set-based detection scheme is used to monitor the process, similar to the previous section. The values of the process disturbances and measurement noise are varied at each time step over this simulation. The values of the monitoring variables at the switching instances over this simulation are illustrated in Figure 6a. The attack is not detected over this simulation, demonstrating that an attacker with sufficient knowledge of the detection scheme may be able to design an attack that is capable of evading detection under a scheduled control mode switch.



**Figure 6.** Values of the monitoring variable for the CSTR process under attack showing (a) no attack detection with scheduled control mode switching (b) attack detection at  $t_d = 235$  under the randomized switching scheme after a control mode switch implemented at  $t_{s_1} = 221$ .

Next, we perform 1000 closed-loop simulations considering the process subject to the smart attack but with randomized (rather than scheduled) control mode switching. Over each simulation, we use the same values of the process disturbances and measurement noise as in the first simulation with the scheduled control mode switch. However, we apply the control mode switches at randomly chosen time steps. Over each simulation, we allow for a maximum of three control mode switches. The first control mode switching instance from the nominal to the attack-sensitive control mode ( $t_{s_1}$ ) is selected as a random integer generated over the interval  $[99, 249]$ . The second switching instance from the nominal to the attack-sensitive control mode ( $t_{s_3}$ ) depends on the first switching instance and is selected as a random integer generated over the interval  $[t_{s_1} + T_c^A + 151, t_{s_1} + T_c^A + 451]$ . Finally, if an attack is not detected over the two previous control mode switches, a third switch from the nominal control mode to the attack-sensitive control mode is implemented at time step ( $t_{s_5}$ ) selected as a random integer generated over the interval  $[t_{s_3} + T_c^A + 151, t_{s_3} + T_c^A + 451]$ . The third switch from the nominal to the attack-sensitive control mode is implemented only if  $t_{s_3} + T_c^A + 151 < 1200$  so that at the end of each simulation, the process is under the nominal control mode.

Over all simulations considering the process under the attack and the randomized control mode switching, the detection of the smart attack occurs after the first control mode switch within a minimum of 4 time steps and a maximum of 40 time steps from the switching instance. Figure 6b illustrates the attack detection over a simulation considering the process under the smart attack with the randomized control mode switching-enabled attack detection method. Over this simulation, the control mode switch is implemented at time step  $t_{s_1} = 221$  leading to attack detection at  $t_d = 235$ . Because an attack is detected after the first control mode switch, no further switches are implemented by the detection scheme. The results demonstrate that from the perspective of enabling attack detection, a randomized control mode switching method may be preferred to a scheduled control mode switching method because an attacker may not be able to design a smart attack that is capable of evading detection.

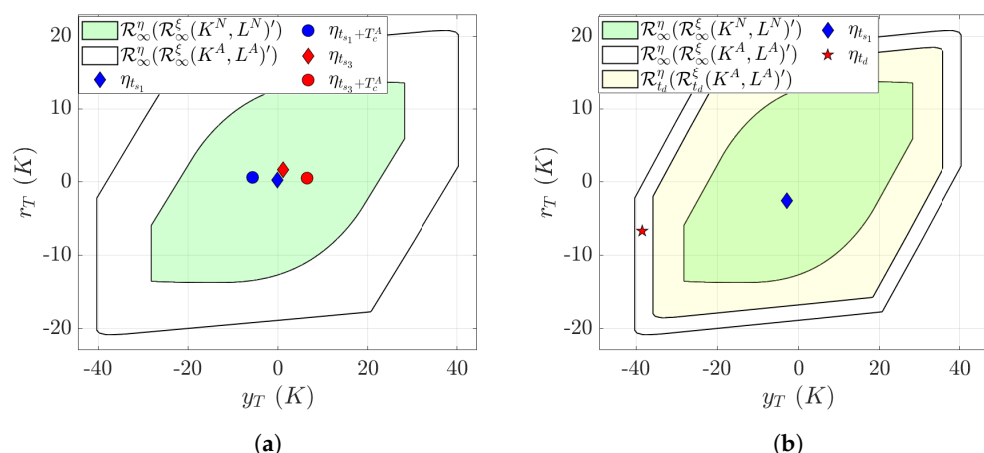
**Remark 5.** Figure 4 is an illustrative example for the operation of the CSTR process under the scheduled switching-enabled detection strategy. An operator may choose to implement a switching-enabled detection strategy on the process that uses other periodic or non-periodic patterns for switching between the nominal and the attack-sensitive modes. Irrespective of the pattern of operation, an attacker with knowledge of the switching schedule may be able to design a detection-evading smart attack. Comparison of attack detection between scheduled and randomized switching-enabled strategies for smart attacks designed to evade detection under scheduled switching with patterns other than the one shown in Figure 4 may require a case-by-case analysis. Nevertheless, we anticipate that operating the CSTR process under the randomized switching-enabled strategy will help preserve the confidentiality of the detection scheme and thus enable attack detection over more simulations than under the scheduled switching-enabled detection strategy.

#### 4.2.2. Application to the Nonlinear CSTR Process

We demonstrate the application of the proposed switching-enabled attack detection method to detect a smart attack on the nonlinear sampled-data model of the CSTR process. In all simulations considered in this section, the CSTR is modeled using its continuous-time nonlinear process model in Equation (21), and the linear control law and the Luengerberger observer are applied with a zero-order hold that considers a sampling interval of  $\Delta = 1 \times 10^{-2}$  h. To solve the differential equations modeling the CSTR, we use the explicit Euler's method with an integration time step of  $1 \times 10^{-4}$  h. We model the process disturbances and measurement noise values as random numbers drawn from a Gaussian distribution. Specifically, we model the disturbance in the feed concentration ( $\Delta C_{A0}$ ) as random variables drawn from a distribution with  $\mathcal{N}(0, 0.0037)$  and the disturbance in the sensor measuring the temperature of the feed to the reactor and the measurement noise in the sensor measuring the temperature of the reactor as random variables drawn from two distinct distributions with  $\mathcal{N}(0, 1.125)$ . The disturbances are drawn from normal distributions with a smaller standard deviation than those used over the simulations in the Section 4.2.1. This is performed to ensure the validity of the reachable and terminal sets computed using the linear process model and enable process monitoring using the detection scheme in Equation (15).

Similar to the previous section, we first simulate the closed-loop process under a smart attack (Figure 5a,b) with a scheduled control switching-enabled attack detection method as in Figure 4. Over this simulation, the closed-loop system is initialized at the origin, and the process disturbances and measurement noise are varied at each sampling instance. As illustrated in Figure 7a, based on the values of the monitoring variables at the switching instances, it can be seen that the smart attack is not detected with the scheduled control mode switching approach. We then perform 1000 simulations of the attacked process with the same process disturbance and measurement noise considered in the simulation with the scheduled switch. Over each simulation, a randomized control mode switching is implemented, where the switching instances are the same as over the corresponding simulations for the linear CSTR process model in Section 4.2.1. Over all simulations, detection of the attack occurs within a minimum of 5 time steps (0.05 h in real time) and a maximum of 44 time steps (0.44 h in real time) from the first switching instance. Figure 7b illustrates the evolution of the monitoring variable over a few time steps of the process under the smart attack, with the control mode switch implemented randomly. Over this simulation, attack detection occurs at time step  $t_d = 231$  (2.31 h in real time), which is 10 time steps after the first switch from the nominal control mode to the attack-sensitive control mode is implemented at time step  $t_{s_1} = 221$  (2.21 h in real time). The reachable set for the attack-free process at the detection time step is contained entirely within the terminal set. The result highlights that, for this simulation, monitoring of the process using a reachable set-based detection scheme may be preferable to monitoring of the process using the terminal set-based detection scheme to aid in the detection of the attack at the earliest time step possible. Following the detection of the attack, the process

is switched back to the nominal control mode and no further alarms are observed in the detection scheme.



**Figure 7.** The values of the monitoring variable for the nonlinear CSTR under a smart attack showing (a) no attack detection with a scheduled control mode switch and (b) attack detection at time step  $t_d = 231$  (2.31 h in real time) under the randomized switching scheme after a control mode switch at time step  $t_{s_1} = 221$  (2.21 h in real time).

## 5. Conclusions

In this work, we presented a cyberattack detection method that utilizes randomized control mode switching to enable the detection of an attack on processes during transient operation. The proposed detection method guarantees no alarms in the detection scheme for the attack-free process with the control mode switching. In developing the detection method, we exploited the interdependence between the control parameters, closed-loop stability of the attacked process, and the ability of a reachable set-based detection scheme to detect the attack. As chemical processes are under prolonged operation at or near their steady states, we proposed a modification of the detection method for application to processes operating under steady-state conditions. Using two illustrative examples, we demonstrated the application of the control mode switching for attack detection with zero false alarms. In the first example, we considered a scalar process under transient operation, while in the second example, we considered a chemical process under steady-state operation. Finally, using simulations of the chemical process example, we demonstrated that a randomized control mode switch may prevent an attacker from learning the switching schedule, thereby preventing them from designing a smart attack that evades detection.

**Author Contributions:** Conceptualization, S.N., M.J.E. and N.H.E.-F.; methodology, S.N., M.J.E. and N.H.E.-F.; formal analysis, S.N., M.J.E. and N.H.E.-F.; software, S.N.; visualization, S.N.; writing—original draft preparation, S.N., M.J.E. and N.H.E.-F.; writing—review and editing, M.J.E. and N.H.E.-F.; project administration, M.J.E. and N.H.E.-F.; supervision, M.J.E. and N.H.E.-F. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by the National Science Foundation grant number CBET-2137281.

**Data Availability Statement:** Dataset available on request from the authors.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

- Setola, R.; Faramondi, L.; Salzano, E.; Cozzani, V. An overview of cyber attack to industrial control system. *Chem. Eng. Trans.* **2019**, *77*, 907–912. [CrossRef]
- Fact Sheet: Biden-Harris Administration Expands Public-Private Cybersecurity Partnership to Chemical Sector. Available online: <https://www.whitehouse.gov/briefing-room/statements-releases/2022/10/26/fact-sheet-biden-harris-administration-expands-public-private-cybersecurity-partnership-to-chemical-sector/> (accessed on 31 January 2024).

3. Tsuchiya, A.; Fraile, F.; Koshijima, I.; Ortiz, A.; Poler, R. Software defined networking firewall for industry 4.0 manufacturing systems. *J. Ind. Eng. Manag.* **2018**, *11*, 318–332. [\[CrossRef\]](#)
4. Albert, R.; Jeong, H.; Barabási, A.L. Error and attack tolerance of complex networks. *Nature* **2000**, *406*, 378–382. [\[CrossRef\]](#) [\[PubMed\]](#)
5. Shang, K.K.; Li, T.C.; Small, M.; Burton, D.; Wang, Y. Link prediction for tree-like networks. *Chaos Interdiscip. J. Nonlinear Sci.* **2019**, *29*, 061103. [\[CrossRef\]](#) [\[PubMed\]](#)
6. Chen, S.; Wu, Z.; Christofides, P.D. A cyber-secure control-detector architecture for nonlinear processes. *AIChE J.* **2020**, *66*, e16907. [\[CrossRef\]](#)
7. Durand, H. Process/Equipment design implications for control system cybersecurity. In Proceedings of the 9th International Conference on Foundations of Computer-Aided Process Design, Copper Mountain, CO, USA, 14–18 July 2019; Volume 47, pp. 263–268. [\[CrossRef\]](#)
8. Wen, H.; Khan, F.; Ahmed, S.; Imtiaz, S.; Pistikopoulos, S. Risk assessment of human-automation conflict under cyberattacks in process systems. *Comput. Chem. Eng.* **2023**, *172*, 108175. [\[CrossRef\]](#)
9. Parker, S.; Wu, Z.; Christofides, P. Cybersecurity in process control, operations, and supply chain. *Comput. Chem. Eng.* **2023**, *171*, 108169. [\[CrossRef\]](#)
10. Rangan, K.K.; Oyama, H.; Durand, H. Integrated cyberattack detection and handling for nonlinear systems with evolving process dynamics under Lyapunov-based economic model predictive control. *Chem. Eng. Res. Des.* **2021**, *170*, 147–179. [\[CrossRef\]](#)
11. Oyama, H.; Rangan, K.K.; Durand, H. Handling of stealthy sensor and actuator cyberattacks on evolving nonlinear process systems. *J. Adv. Manuf. Process.* **2021**, *3*, e10099. [\[CrossRef\]](#)
12. Oyama, H.; Messina, D.; Rangan, K.K.; Durand, H. Lyapunov-based economic model predictive control for detecting and handling actuator and simultaneous sensor/actuator cyberattacks on process control systems. *Front. Chem. Eng.* **2022**, *4*, 810129. [\[CrossRef\]](#)
13. Narasimhan, S.; El-Farra, N.H.; Ellis, M.J. A control-switching approach for cyberattack detection in process systems with minimal false alarms. *AIChE J.* **2022**, *68*, e17875. [\[CrossRef\]](#)
14. Narasimhan, S.; El-Farra, N.H.; Ellis, M.J. A reachable set-based scheme for the detection of false data injection cyberattacks on dynamic processes. *Dig. Chem. Eng.* **2023**, *7*, 100100. [\[CrossRef\]](#)
15. Narasimhan, S.; El-Farra, N.H.; Ellis, M.J. Active multiplicative cyberattack detection utilizing controller switching for process systems. *J. Process Control* **2022**, *116*, 64–79. [\[CrossRef\]](#)
16. Narasimhan, S.; El-Farra, N.H.; Ellis, M.J. Detectability-based controller design screening for processes under multiplicative cyberattacks. *AIChE J.* **2022**, *68*, e17430. [\[CrossRef\]](#)
17. Zedan, A.; El-Farra, N.H. A machine-learning approach for identification and mitigation of cyberattacks in networked process control systems. *Chem. Eng. Res. Des.* **2021**, *176*, 102–115. [\[CrossRef\]](#)
18. Mavridis, C.N.; Kanellopoulos, A.; Vamvoudakis, K.G.; Baras, J.S.; Johansson, K.H. Attack identification for cyber-physical security in dynamic games under cognitive hierarchy. In Proceedings of the IFAC World Congress, Yokohama, Japan, 9–14 July 2023; In Press.
19. Cómbita, L.F.; Quijano, N.; Cárdenas, Á.A. On the stability of cyber-physical control systems with sensor multiplicative attacks. *IEEE Access* **2022**, *10*, 39716–39728. [\[CrossRef\]](#)
20. Liu, H.; Mo, Y.; Johansson, K.H. Active detection against replay attack: A survey on watermark design for cyber-physical systems. In *Lecture Notes in Control and Information Sciences*; Springer International Publishing: Cham, Switzerland, 2021; pp. 145–171. [\[CrossRef\]](#)
21. Mo, Y.; Sinopoli, B. Secure control against replay attacks. In Proceedings of the 47th Annual Allerton Conference on Communication, Control, and Computing, Monticello, IL, USA, 30 September–2 October 2009; pp. 911–918. [\[CrossRef\]](#)
22. Ko, W.H.; Satchidanandan, B.; Kumar, P.R. Dynamic watermarking-based defense of transportation cyber-physical systems. *ACM Trans. Cyber-Phys. Syst.* **2019**, *4*, 1–21. [\[CrossRef\]](#)
23. Oyama, H.; Messina, D.; Rangan, K.K.; Akkarakaran, F.L.; Nieman, K.; Durand, H.; Tyrrell, K.; Hinzman, K.; Williamson, M. Development of directed randomization for discussing a minimal security architecture. *Dig. Chem. Eng.* **2023**, *6*, 100065. [\[CrossRef\]](#)
24. Huang, T.; Satchidanandan, B.; Kumar, P.R.; Xie, L. An online detection framework for cyber attacks on automatic generation control. *IEEE Trans. Power Syst.* **2018**, *33*, 6816–6827. [\[CrossRef\]](#)
25. Na, G.; Eun, Y. A multiplicative coordinated stealthy attack and its detection for cyber physical systems. In Proceedings of the IEEE Conference on Control Technology and Applications, Copenhagen, Denmark, 21–24 August 2018; pp. 1698–1703. [\[CrossRef\]](#)
26. Gönen, S.; Sayan, H.H.; Yılmaz, E.N.; Üstünsoy, F.; Karacayılmaz, G. False data injection attacks and the insider threat in smart systems. *Comput. Secur.* **2020**, *97*, 101955. [\[CrossRef\]](#)
27. Reda, H.T.; Anwar, A.; Mahmood, M. Comprehensive survey and taxonomies of false data injection attacks in smart grids: Attack models, targets, and impacts. *Renew. Sust. Energ. Rev.* **2022**, *163*, 112423. [\[CrossRef\]](#)
28. Mo, Y.; Weerakkody, S.; Sinopoli, B. Physical authentication of control systems: Designing watermarked control inputs to detect counterfeit sensor outputs. *IEEE Contr. Syst. Mag.* **2015**, *35*, 93–109. [\[CrossRef\]](#)
29. Durand, H. A nonlinear systems framework for cyberattack prevention for chemical process control systems. *Mathematics* **2018**, *6*, 169. [\[CrossRef\]](#)



30. Guibas, L.J.; Nguyen, A.T.; Zhang, L. Zonotopes as bounding volumes. In Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, Baltimore, MD, USA, 12–14 January 2003; Volume 3, pp. 803–812.
31. Althoff, M.; Frehse, G.; Girard, A. Set propagation techniques for reachability analysis. *Annu. Rev. Control Robot. Auton. Syst.* **2021**, *4*, 369–395. [[CrossRef](#)]
32. Raković, S.V.; Kerrigan, E.C.; Kouramas, K.I.; Mayne, D.Q. Invariant approximations of the minimal robust positively invariant set. *IEEE Trans. Automat. Contr.* **2005**, *50*, 406–410. [[CrossRef](#)]
33. Kuntsevich, V.M.; Pshenichnyi, B.N. Minimal invariant sets of dynamic systems with bounded disturbances. *Cybern. Syst. Anal.* **1996**, *32*, 58–64. [[CrossRef](#)]
34. Lin, H.; Antsaklis, P.J. Stability and stabilizability of switched linear systems: A survey of recent results. *IEEE Trans. Automat. Contr.* **2009**, *54*, 308–322. [[CrossRef](#)]
35. Attar, M.; Lucia, W. Data-driven robust backward reachable sets for set-theoretic model predictive control. *IEEE Contr. Syst. Lett.* **2023**, *7*, 2395–2310. [[CrossRef](#)]
36. Savković, B. Low complexity parameterized approximations of reachable sets for LTI systems. In Proceedings of the 2009 IEEE International Conference on Control and Automation, Christchurch, New Zealand, 9–11 December 2009; pp. 960–965. [[CrossRef](#)]
37. Girard, A.; Le, G.C.; Maler, O. Efficient computation of reachable sets of linear time-invariant systems with inputs. In Proceedings of the 9th International Workshop on Hybrid Systems: Computation and Control, Santa Barbara, CA, USA, 29–31 March 2006; pp. 257–271. [[CrossRef](#)]
38. Herceg, M.; Kvasnica, M.; Jones, C.N.; Morari, M. Multi-Parametric Toolbox 3.0. In Proceedings of the European Control Conference, Zürich, Switzerland, 17–19 July 2013; pp. 502–510. Available online: <http://control.ee.ethz.ch/~mpt> (accessed on 20 January 2024).
39. Althoff, M. An introduction to CORA 2015. *ARCH@ CPSweek* **2015**, *34*, 120–151.
40. Raković, S.V.; Kouramas, K.I. The minimal robust positively invariant set for linear discrete time systems: Approximation methods and control Applications. In Proceedings of the 45th IEEE Conference on Decision and Control, San Diego, CA, USA, 13–15 December 2006; pp. 4562–4567. [[CrossRef](#)]

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.