

Attack Models and Scenarios for Networked Control Systems

André Teixeira, Daniel Pérez, Henrik Sandberg, Karl H. Johansson
ACCESS Linnaeus Center
Automatic Control Laboratory
KTH - Royal Institute of Technology
Osquldas väg 10
SE-10044 Stockholm, Sweden
andretei, danielph, hsan, kallej@kth.se

ABSTRACT

Cyber-secure networked control is modeled, analyzed, and experimentally illustrated in this paper. An attack space defined by the adversary's system knowledge, disclosure, and disruption resources is introduced. Adversaries constrained by these resources are modeled for a networked control system architecture. It is shown that attack scenarios corresponding to replay, zero dynamics, and bias injection attacks can be analyzed using this framework. An experimental setup based on a quadruple-tank process controlled over a wireless network is used to illustrate the attack scenarios, their consequences, and potential counter-measures.

Categories and Subject Descriptors

K.6.5 [Management of Computing and Information Systems]: Security and Protection—*unauthorized access*;
C.3 [Special-Purpose and Application-Based Systems]: Process control systems

Keywords

Cyber-physical systems, security, attack space, secure control systems

1. INTRODUCTION

Safe and reliable operation of infrastructures is of major societal importance. These systems need to be engineered in such a way so that they can be continuously monitored, coordinated, and controlled despite a variety of potential system disturbances. Given the strict operating requirements and system complexity, such systems are operated through IT infrastructures enabling the timely data flow between digital controllers, sensors, and actuators. However, the use of non-proprietary communication networks and heterogeneous IT components has made these cyber-physical systems vulnerable to cyber threats. One such example are the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

HiCoNS'12, April 17–18, 2012, Beijing, China.

Copyright 2012 ACM 978-1-4503-1263-9/12/04 ...\$10.00.

power transmission networks operated through Supervisory Control and Data Acquisition (SCADA) systems. The measurement and control data in these systems are commonly transmitted through unprotected channels, leaving the system vulnerable to several threats [7]. In fact cyber attacks on power networks operated by SCADA systems have been reported in the media [8].

There exists a vast literature on computer security focusing on three main properties of data and IT services, namely confidentiality, integrity, and availability [3]. Confidentiality relates to the non-disclosure of data by unauthorized parties. Integrity on the other hand concerns the trustworthiness of data, meaning there is no unauthorized change of the data contents or properties, while availability means that timely access to the data or system functionalities is ensured. Unlike other IT systems where cyber-security mainly involves the protection of data, cyber attacks on networked control systems may influence the physical processes through the communication infrastructure due to feedback loops. Therefore networked control system security needs to consider the existing threats at both the cyber and physical layers. These threats can be captured in the attack space illustrated in Figure 1, which also depicts several attack scenarios described in this work. For instance, two typical examples of cyber attacks considered in IT security can be found in Figure 1, the eavesdropping attack and the Denial-of-Service (DoS) attack.

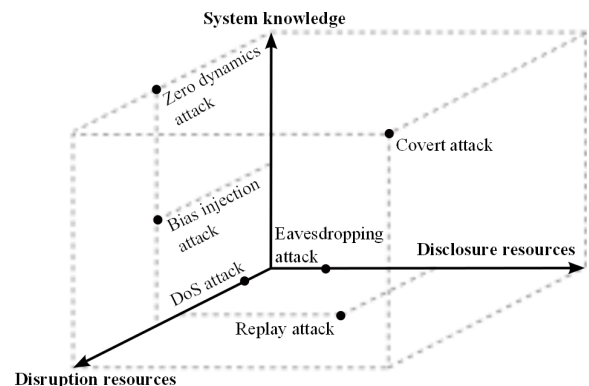


Figure 1: The cyber-physical attack space.

We propose three dimensions for the attack space: the adversary's *a priori* system model knowledge and his disclo-

sure and disruption resources. The *a priori* system knowledge can be used by the attacker to construct more complex attacks, possibly harder to detect and with more severe consequences. Similarly, the disclosure resources enable the attacker to obtain sensitive information about the system during the attack by violating the data confidentiality. Note that disclosure resources cannot be used to disrupt the system operation, which is the case of the eavesdropping attack illustrated in Figure 1. On the other hand, disruption resources can be used to affect the system operation, which happens for instance when data integrity or availability properties are violated. One such example is the DoS attack illustrated in Figure 1, where the data required for correctly operating the system is made unavailable.

Control theory has contributed with frameworks to handle model uncertainties and disturbances [20] as well as fault diagnosis and mitigation [5, 10], which can then be used to detect and attenuate the consequences of cyber attacks on networked control systems. Some of these tools are therefore considered as part of the networked control system and will be used to analyze the consequences of cyber attacks.

1.1 Related Work

Cyber attacks on control systems compromising measurement and actuator data integrity and availability have been considered in [4], where the authors modeled their effects on the physical dynamics. Availability attacks have been further analyzed in [1, 9] for resource constrained attackers with full-state information. Particularly, this work considered DoS attacks in which the attacker could jam the communication channels and prevent measurement and actuator data from reaching its destination, rendering the data unavailable. A particular instance of the DoS attack in which the attacker does not have any *a priori* system knowledge [1] is represented in the attack space in Figure 1.

Deception attacks compromising integrity have recently received more attention. A particular kind of deception attacks, i.e. replay attacks on the sensor measurements, has been analyzed in [14]. The authors considered the case where all the existing sensors were attacked and proposed suitable counter-measures to detect the attack. In this attack scenario the attacker does not have any system knowledge but is able to access and corrupt the sensor data, thus having disclosure and disruptive resources, as depicted in Figure 1.

Another class of deception attacks, false-data injection attacks, has also been studied in recent work. For instance, in the case of power networks, an attacker with perfect model knowledge has been initially considered in [13]. The work in [12] considered stealthy attacks with limited resources and proposed improved detection methods, while [16] analyzed the minimum number of sensors required for stealthy attacks, based on which measurement security metrics were proposed. The consequences of these attacks have also been analyzed in [18, 19]. The models used are static, hence these attack scenarios are closest to the bias injection attack shown in Figure 1.

Data injection attacks on dynamic control systems were also considered. In [17] the author characterizes the set of attack policies for covert (undetectable) false-data injection attacks with detailed model knowledge and full access to all sensor and actuator channels, while [15] described the set of undetectable false-data injection attacks for omniscient

attackers with full-state information, but possibly compromising only a subset of the existing sensors and actuators. In these attack scenarios confidentiality was also violated, as the attacker had access to either measurement and actuator data or full-state information. These attacks are therefore placed close to the boundaries of the attack space, as illustrated in Figure 1 for the covert attack, while the framework in [15] addresses attacks on the top plane where full model knowledge is considered.

1.2 Contributions and Outline

Most of the recent work on cyber-security of control systems has considered scenarios where the attacker has access to a large set of resources and knowledge, thus being placed close to the boundaries of the attack space in Figure 1. Therefore a large part of the attack space has not been addressed. In particular, the class of detectable attacks that do not trigger conventional alarms has yet to be covered in depth.

In this paper we consider a typical control architecture for the networked control system under both cyber and physical attacks. Given this architecture, a generic adversary model applicable to several attack scenarios is discussed and the attack resources are mapped to the corresponding dimensions of the attack space. Three stealthy attack scenarios are discussed in more detail to better illustrate the proposed adversary model and the concept of attack space.

One of the attack scenarios analyzed corresponds to a particular type of detectable attack, the bias injection attack. Although this attack may be detected, it requires limited model knowledge and no information about the system state. Stealthiness conditions are provided, as well as a methodology to assess the attack impact on the physical state of the system.

The attack scenarios analyzed in the paper have been staged at our testbed for security of control systems. The testbed architecture and results from the staged attacks are presented and discussed.

The outline of the paper is as follows. The system architecture and model are described in Section 2, while Section 3 contains the adversary model and a detailed description of the attack resources on each dimension of the attack space. The framework introduced in the previous sections is then illustrated for three particular attack scenarios in Section 4. The results of the experiments for each attack scenario in a secure control systems testbed are presented and discussed in Section 5, followed by conclusions in Section 6.

2. NETWORKED CONTROL SYSTEM

In this section we describe the networked control system structure, where we consider three main components: the physical plant and communication network, the feedback controller, and the anomaly detector.

2.1 Physical Plant and Communication Network

The physical plant is modeled in a discrete-time state-space form,

$$\mathcal{P} : \begin{cases} x_{k+1} = Ax_k + B\tilde{u}_k + Gw_k + Ff_k \\ y_k = Cx_k + v_k \end{cases}, \quad (1)$$

where $x_k \in \mathbb{R}^n$ is the state variable, $\tilde{u}_k \in \mathbb{R}^q$ the control actions applied to the process, $y_k \in \mathbb{R}^p$ the measurements

from the sensors at the sampling instant $k \in \mathbb{Z}$, and $f_k \in \mathbb{R}^d$ is the unknown signal representing the effects of anomalies, usually denoted as fault signal in the fault diagnosis literature [6]. The process and measurement noise, $w_k \in \mathbb{R}^n$ and $v_k \in \mathbb{R}^p$, represent the discrepancies between the model and the real process, due to unmodeled dynamics or disturbances, for instance, and we assume their means are respectively bounded by δ_w and δ_v , i.e. $\bar{w} = \|\mathbb{E}\{w_k\}\| \leq \delta_w$ and $\bar{v} = \|\mathbb{E}\{v_k\}\| \leq \delta_v$.

The physical plant operation is supported by a communication network through which the sensor measurements and actuator data are transmitted, which at the plant side correspond to y_k and \tilde{u}_k , respectively. At the controller side we denote the sensor and actuator data by $\tilde{y}_k \in \mathbb{R}^p$ and $u_k \in \mathbb{R}^q$ respectively. Since the communication network may be unreliable, the data exchanged between the plant and the controller may be altered, resulting in discrepancies in the data at the plant and controller ends. In this paper we do not consider the usual communication network effects such as packet losses and delays. Instead we focus on data corruption due to malicious cyber attacks, as described in Section 3. Therefore the communication network is supposed to be reliable, not affecting the data flowing through it.

Given the physical plant model (1) and assuming an ideal communication network, the networked control system is said to have a *nominal behavior* if $f_k = 0$, $\tilde{u}_k = u_k$, and $\tilde{y}_k = y_k$. The absence of either one of these condition results in an abnormal behavior of the system.

2.2 Feedback Controller

In order to comply with performance requirements in the presence of the unknown process and measurement noises, we consider that the physical plant is controlled by an appropriate linear time-invariant feedback controller [20]. The output feedback controller can be written in a state-space form as

$$\mathcal{F} : \begin{cases} z_{k+1} = A_c z_k + B_c \tilde{y}_k \\ u_k = C_c z_k + D_c \tilde{y}_k \end{cases} \quad (2)$$

where the states of the controller, $z_k \in \mathbb{R}^m$, may include the process state and tracking error estimates. Given the plant and communication network models, the controller is supposed to be designed so that acceptable performance is achieved under nominal behavior.

2.3 Anomaly Detector

In this section we consider the anomaly detector that monitors the system to detect possible anomalies, i.e. deviations from the nominal behavior. We consider that the anomaly detector is collocated with the controller, therefore it only has access to \tilde{y}_k and u_k to evaluate the behavior of the plant.

Several approaches to detecting malfunctions in control systems are available in the fault diagnosis literature [6, 10]. Here we consider the following observer-based Fault Detection Filter

$$\mathcal{D} : \begin{cases} \hat{x}_{k|k} = A \hat{x}_{k-1|k-1} + B u_{k-1} + K(\tilde{y}_k - \hat{y}_{k|k-1}) \\ r_k = V(\tilde{y}_k - \hat{y}_{k|k}) \end{cases}, \quad (3)$$

where $\hat{x}_{k|k} \in \mathbb{R}^n$ is the state estimate given measurements up until time k and $r_k \in \mathbb{R}^p$ the residue, which is evaluated in order to detect and locate existing anomalies.

The anomaly detector is designed so that

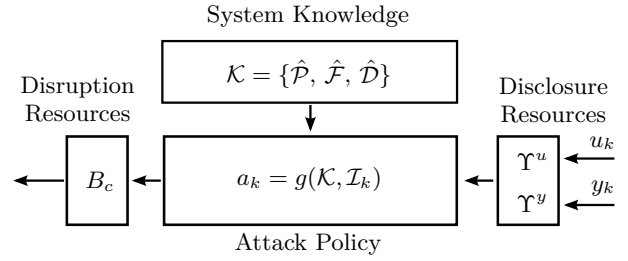


Figure 2: Adversary model for a point in the attack space in Figure 1.

1. under nominal behavior of the system (i.e., $f_k = 0$, $u_k = \tilde{u}_k$, $y_k = \tilde{y}_k$), the expected value of the residue converges asymptotically to a neighborhood of zero, i.e., $\lim_{k \rightarrow \infty} \|\mathbb{E}\{r_k\}\| \leq \delta_r$, with $\delta_r \in \mathbb{R}^+$;
2. the residue is sensitive to the anomalies ($f_k \neq 0$).

An alarm is triggered if the residue meets

$$\|r_k\| \geq \delta_r + \delta_\alpha, \quad (4)$$

where $\delta_\alpha \in \mathbb{R}^+$ is chosen so that the false alarm rate does not exceed a given $\alpha \in [0, 1]$.

3. ADVERSARY MODELS

The adversary model considered in this paper is illustrated in Figure 2 and is composed of an attack policy and the adversary resources i.e., the system model knowledge, the disclosure resources, and the disruption resources. The attack policy is described by

$$a_k = g(\mathcal{K}, \mathcal{I}_k). \quad (5)$$

Each of the attack policy components can be mapped to a specific axis of the attack space in Figure 1: $\mathcal{K} = \{\hat{\mathcal{P}}, \hat{\mathcal{F}}, \hat{\mathcal{D}}\}$ is the *a priori* system knowledge possessed by the attacker, \mathcal{I}_k corresponds to the set of sensor and actuator data available to the attacker at time k , thus being mapped to the disclosure resources, while a_k is the attack vector at time k that may affect the system behavior using the disruption resources captured by B_c .

In this section we describe the networked control system under attack with respect to the attack vector a_k . Then we detail the adversary's system knowledge, the disclosure resources, and the disruption resources. Models of the attack vector a_k for particular disruption resources are also given.

3.1 Networked Control System under Attack

The system components under attack are now characterized for the attack vector a_k . Considering the plant and controller states to be stacked as $\eta_k = [x_k^\top \ z_k^\top]^\top$, the dynamics of the closed-loop system composed by \mathcal{P} and \mathcal{F} under the effect of a_k can be written as

$$\begin{aligned} \eta_{k+1} &= \mathbf{A}_c \eta_k + \mathbf{B}_c a_k + \mathbf{G}_c \begin{bmatrix} w_k \\ v_k \end{bmatrix} \\ \tilde{y}_k &= \mathbf{C}_c \eta_k + \mathbf{D}_c a_k + \mathbf{H}_c \begin{bmatrix} w_k \\ v_k \end{bmatrix}, \end{aligned} \quad (6)$$

where the system matrices are

$$\mathbf{A}_c = \begin{bmatrix} A + BD_cC & BC_c \\ B_cC & A_c \end{bmatrix}, \quad \mathbf{G}_c = \begin{bmatrix} G & BD_c \\ 0 & B_c \end{bmatrix},$$

$$\mathbf{C}_c = [C \ 0], \quad \mathbf{H}_c = [0 \ I],$$

and \mathbf{B}_c and \mathbf{D}_c capture the way in which the attack vector a_k affects the plant and controller. These matrices are characterized for some attack scenarios in Section 3.4.

Similarly, using \mathcal{P} and \mathcal{D} as in (1) and (3), respectively, the anomaly detector error dynamics under attack are described by

$$\xi_{k|k} = \mathbf{A}_e \xi_{k-1|k-1} + \mathbf{B}_e a_{k-1} + \mathbf{G}_e \begin{bmatrix} w_{k-1} \\ v_k \end{bmatrix}$$

$$r_k = \mathbf{C}_e \xi_{k-1|k-1} + \mathbf{D}_e a_{k-1} + \mathbf{H}_e \begin{bmatrix} w_{k-1} \\ v_k \end{bmatrix}, \quad (7)$$

where $\xi_{k|k} \in \mathbb{R}^n$ is the estimation error and

$$\mathbf{A}_e = (I - KC)A, \quad \mathbf{G}_e = [(I - KC)G \ -K],$$

$$\mathbf{C}_e = VC(I - KC)A, \quad \mathbf{H}_e = [VC(I - KC)G \ V(I - CK)].$$

The matrices \mathbf{B}_e and \mathbf{D}_e are specific to the available disruptive resources and are characterized in Section 3.4.

3.2 System Knowledge

The amount of *a priori* knowledge regarding the control system is a core component of the adversary model, as it may be used, for instance, to render the attack undetectable. In general, we may consider that the adversary approximately knows the model of the plant ($\hat{\mathcal{P}}$) and the algorithms used in the feedback controller ($\hat{\mathcal{F}}$) and the anomaly detector ($\hat{\mathcal{D}}$), thus denoting the adversary knowledge by $\mathcal{K} = \{\hat{\mathcal{P}}, \hat{\mathcal{F}}, \hat{\mathcal{D}}\}$. Figure 1 illustrates several types of attack scenarios with different amounts of required system knowledge. In particular, note that the replay attacks do not need much knowledge of the system components.

3.3 Disclosure Resources

The disclosure resources enable the attacker to gather sequences of data from the calculated control actions u_k and the real measurements y_k through disclosure attacks. Denote $\mathcal{R}_C^u \subseteq \{1, \dots, q\}$ and $\mathcal{R}_C^y \subseteq \{1, \dots, p\}$ as the disclosure resources, i.e. set of actuator and sensor channels that can be accessed during disclosure attacks, and let \mathcal{I}_k be the control and measurement data sequence gathered by the attacker from time k_0 to k . The disclosure attacks can then be modeled as

$$\mathcal{I}_k := \mathcal{I}_{k-1} \cup \begin{bmatrix} \Upsilon^u & 0 \\ 0 & \Upsilon^y \end{bmatrix} \begin{bmatrix} u_k \\ y_k \end{bmatrix}, \quad (8)$$

where $\Upsilon^u \in \mathbb{B}^{|\mathcal{R}_C^u| \times q}$ and $\Upsilon^y \in \mathbb{B}^{|\mathcal{R}_C^y| \times p}$ are the binary incidence matrices mapping the data channels to the corresponding data gathered by the attacker and $\mathcal{I}_{k_0} = \emptyset$.

As seen in the above description of disclosure attacks, the physical dynamics of the system are not affected by these type of attacks. Instead, these attacks gather intelligence that may enable more complex attacks, such as the replay attacks depicted in Figure 1.

3.4 Disruption Resources

As seen in (6) and (7), disruption resources are related to the attack vector a_k and may be used to affect the several components of the system. The way a particular attack disturbs the system operation depends not only on the respective resources, but also on the nature of the attack. For instance, a physical attack directly perturbs the system dynamics, whereas a cyber attack disturbs the system through the cyber-physical couplings. To better illustrate this discussion we now consider physical, data deception, and data DoS attacks.

3.4.1 Physical Attack

Physical attacks may occur in control systems, often in conjunction with cyber attacks. For instance, in [2] water was pumped out of an irrigation system while the water level measurements were corrupted so that the attack remained stealthy. Since physical attacks are similar to the fault signals f_k in (1), in the following sections we consider f_k to be the physical attack modifying the plant dynamics as

$$x_{k+1} = Ax_k + B\tilde{u}_k + Gw_k + Ff_k$$

$$y_k = Cx_k.$$

Considering $a_k = f_k$, the resulting system dynamics are described by (6) and (7) with

$$\mathbf{B}_c = \begin{bmatrix} F \\ 0 \end{bmatrix}, \quad \mathbf{D}_c = 0,$$

$$\mathbf{B}_e = (I - KC)F, \quad \mathbf{D}_e = VC(I - KC)F.$$

Note that the disruption resources in this attack are captured in the matrix F .

3.4.2 Data Deception Attack

The deception attacks modify the control actions u_k and sensor measurements y_k from their calculated or real values to the corrupted signals \tilde{u}_k and \tilde{y}_k , respectively. Denoting $\mathcal{R}_I^u \subseteq \{1, \dots, q\}$ and $\mathcal{R}_I^y \subseteq \{1, \dots, p\}$ as the deception resources, i.e. set of actuator and sensor channels that can be affected, the deception attacks are modeled as

$$\tilde{u}_k := u_k + \Gamma^u b_k^u$$

$$\tilde{y}_k := y_k + \Gamma^y b_k^y \quad (9)$$

where the signals $b_k^u \in \mathbb{R}^{|\mathcal{R}_I^u|}$ and $b_k^y \in \mathbb{R}^{|\mathcal{R}_I^y|}$ represent the data corruption and $\Gamma^u \in \mathbb{B}^{q \times |\mathcal{R}_I^u|}$ and $\Gamma^y \in \mathbb{B}^{p \times |\mathcal{R}_I^y|}$ ($\mathbb{B} := \{0, 1\}$) are the binary incidence matrices mapping the data corruption to the respective data channels. The matrices Γ^u and Γ^y indicate which data channels can be accessed by the attacker and are therefore directly related to the attacker resources in deception attacks.

Defining $a_k = [b_k^{u\top} \ b_{k+1}^{y\top} \ b_k^{y\top}]^\top$, the system dynamics are given by (6) and (7) with

$$\mathbf{B}_c = \begin{bmatrix} B\Gamma^u & 0 & BD_c\Gamma^y \\ 0 & 0 & B_c\Gamma^y \end{bmatrix}, \quad \mathbf{D}_c = [0 \ 0 \ \Gamma^y],$$

$$\mathbf{B}_e = [(I - KC)B\Gamma^u \ -K\Gamma^y \ 0],$$

$$\mathbf{D}_e = [VC(I - KC)B\Gamma^u \ V(I - CK)\Gamma^y \ 0].$$

Note that deception attacks do not possess any disclosure capabilities, as depicted in Figure 1 for examples of deception attacks such as the bias injection attack.

3.4.3 Data Denial-of-Service Attack

The DoS attacks prevent the actuator and sensor data from reaching their respective destinations and should therefore be modeled as the absence of data, for instance $u_k = \emptyset$ if all the actuator data was jammed. However such a model would not fit the framework in (6) and (7) where a_k is assumed to be a real valued vector. Hence we consider instead one of the typical mechanisms used by digital controllers to deal with the absence of data, in which the absent data is considered to be zero. Denoting $\mathcal{R}_A^u \subseteq \{1, \dots, q\}$ and $\mathcal{R}_A^y \subseteq \{1, \dots, p\}$ as the set of actuator and sensor channels that can be jammed, we can model DoS attacks as deception attacks in (9) with

$$\begin{aligned} b_k^u &:= -S_k^u \Gamma^{u\top} u_k \\ b_k^y &:= -S_k^y \Gamma^{y\top} y_k \end{aligned} \quad (10)$$

where $S_k^u \in \mathbb{B}^{|\mathcal{R}_A^u| \times |\mathcal{R}_A^u|}$ and $S_k^y \in \mathbb{B}^{|\mathcal{R}_A^y| \times |\mathcal{R}_A^y|}$ are boolean diagonal matrices where the i -th diagonal entry indicates whether a DoS attack is performed ($[S_k^{(\cdot)}]_{ii} = 1$) or not ($[S_k^{(\cdot)}]_{ii} = 0$) on the corresponding channel. Therefore DoS attacks on the data are a type of disruptive attacks, as depicted in Figure 1.

4. ATTACK SCENARIOS

In this section we discuss the general goal of an attacker and likely choices of the attack policy $g(\cdot, \cdot)$. In particular we consider three attack scenarios with stealthiness constraints under the framework introduced in the previous sections. For each scenario we comment on the attacker's capabilities along each dimension of the attack space in Figure 1 and formulate the corresponding stealthy attack policy. These scenarios are illustrated by experiments on a process control testbed in Section 5.

4.1 Attack Goals and Constraints

The attack scenarios need to also include the intent of the attacker, namely the attack goals and constraints shaping the attack policy. The attack goals can be stated in terms of the attack impact on the system operation, while the constraints may be related to the attack detectability. In this paper we focus on the latter and consider stealthy attacks. Furthermore, we consider the disruptive attack component consists of only physical and data deception attacks, and thus we consider the attack vector $a_k = [f_k^\top \ b_k^{u\top} \ b_{k+1}^{y\top} \ b_k^{y\top}]^\top$.

Given the anomaly detector described in Section 2 and denoting $\mathcal{A}_{k_0}^{k_f} = \{a_{k_0}, \dots, a_{k_f}\}$ as the attack signal, the set of stealthy attacks are defined as follows.

DEFINITION 1. *The attack signal $\mathcal{A}_{k_0}^{k_f}$ is stealthy if $\|r_k\| < \delta_r + \delta_\alpha \forall k \geq k_0$.*

Note that the above definition is dependent on the initial state of the system at k_0 , as well as the noise terms w_k and v_k .

Since the closed-loop system (6) and the anomaly detector (7) under physical and data deception attacks are linear systems, each of these systems can be separated in two components, the nominal component with $a_k = 0 \forall k$ and the following systems

$$\begin{aligned} \eta_{k+1} &= \mathbf{A}_c \eta_k + \mathbf{B}_c a_k \\ \tilde{y}_k^a &= \mathbf{C}_c \eta_k + \mathbf{D}_c a_k \end{aligned} \quad (11)$$

and

$$\begin{aligned} \xi_{k|k} &= \mathbf{A}_e \xi_{k-1|k-1} + \mathbf{B}_e a_{k-1} \\ r_k^a &= \mathbf{C}_e \xi_{k-1|k-1} + \mathbf{D}_e a_{k-1}, \end{aligned} \quad (12)$$

with $\eta_0 = \xi_{0|0} = 0$.

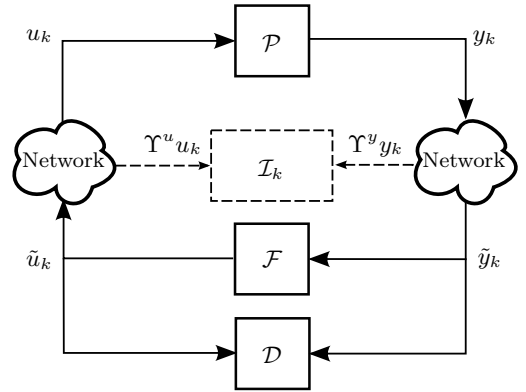
Assuming the system to be in nominal behavior before the attack, using the triangle inequality and linearity property we have $\|r_k^a\| \leq \delta_\alpha \Rightarrow \|r_k\| \leq \delta_r + \delta_\alpha$, leading to the following definition:

DEFINITION 2. *The attack signal $\mathcal{A}_{k_0}^{k_f}$ is α -stealthy with respect to \mathcal{D} if $\|r_k^a\| < \delta_\alpha \forall k \geq k_0$.*

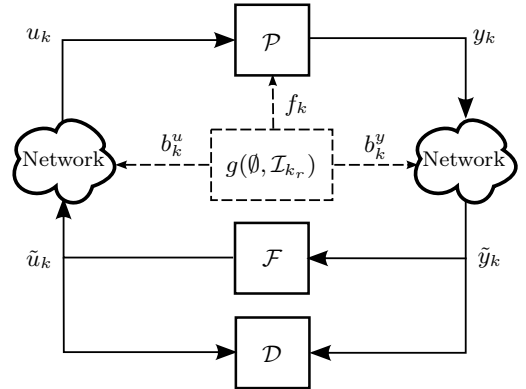
Albeit more conservative than Definition 1, this definition only depends on the attack signals $\mathcal{A}_{k_0}^{k_f}$. Similarly, the impact of attacks on the closed-loop system can also be analyzed by looking at the linear system (11).

4.2 Replay Attack

In replay attacks the adversary first performs a disclosure attack from $k = k_0$ until k_r , gathering sequences of data \mathcal{I}_{k_r} , and then begins replaying the recorded data at time $k = k_r + 1$ until the end of the attack at $k = k_f$, as illustrated in Figure 3. In the scenario considered here, the attacker is also able to perform a physical attack while replaying the recorded data.



(a) Phase I of the replay attack (13).



(b) Phase II of the replay attack (14).

Figure 3: Schematic of the replay attack.

Attack policy

Similar to the work in [14], assuming $\mathcal{R}_C^{(\cdot)} = \mathcal{R}_I^{(\cdot)}$, meaning that the attacker can corrupt the digital channels from which

the data sequences are gathered, the replay attack policy can be described as

$$\begin{aligned} \text{Phase I: } \quad & a_k = 0 \\ & \mathcal{I}_k = \mathcal{I}_{k-1} \cup \begin{bmatrix} \Upsilon^u & 0 \\ 0 & \Upsilon^y \end{bmatrix} \begin{bmatrix} u_k \\ y_k \end{bmatrix}, \end{aligned} \quad (13)$$

with $k_0 \leq k \leq k_r$ and $\mathcal{I}_{k_0} = \emptyset$ and

$$\begin{aligned} \text{Phase II: } \quad & a_k = \begin{bmatrix} g_f(\mathcal{K}, \mathcal{I}_{k_r}) \\ \Upsilon^u(u_{k-T} - u_k) \\ \Upsilon^y(y_{k+1-T} - y_{k+1}) \\ \Upsilon^y(y_{k-T} - y_k) \end{bmatrix} \\ & \mathcal{I}_k = \mathcal{I}_{k-1}, \end{aligned} \quad (14)$$

where $T = k_r - 1 + k_0$ and $k_r + 1 \leq k \leq k_f$. An interesting instance of this attack scenario consists of applying a pre-defined physical attack to the plant, while using replay attacks to render the attack stealthy. In this case the physical attack signal f_k corresponds to an open-loop signal, $f_k = g_f(k)$.

Disclosure resources

The disclosure capabilities required to stage this attack correspond to the data channels that can be eavesdropped by the attacks, namely \mathcal{R}_C^u and \mathcal{R}_C^y .

Disruption resources

In this case the deception capabilities correspond to the data channels that the attacker can tamper, \mathcal{R}_I^u and \mathcal{R}_I^y . In particular, for replay attacks the attacker can only tamper data channels from which data has been previously recorded, i.e. $\mathcal{R}_I^u \subseteq \mathcal{R}_C^u$ and $\mathcal{R}_I^y \subseteq \mathcal{R}_C^y$.

Direct disruption of the physical system through the signal f_k depends on direct access to the physical system, modeled by the matrix F in (1).

System knowledge

Note that no *a priori* knowledge on the system model is needed for the cyber component of the attack, namely the data disclosure and deception attack, as seen in the attack policy (13) and (14). As for the physical attack, f_k , the required knowledge is scenario dependent. In the scenario considered in the experiments described in Section 5, this component was modeled as an open-loop signal, $f_k = g_f(k)$.

Stealthiness constraints

The work in [14] provided conditions under which replay attacks with access to all measurement data channels are stealthy. However, these attacks are not guaranteed to be stealthy when only a subset of the data channels is attacked. In this case, the stealthiness constraint may require additional knowledge of the system model. For instance, the experiment presented in Section 5 required knowledge of the physical system structure, so that f_k only excited the attacked measurements.

4.3 Zero Dynamics Attack

Recalling that for attacks with only physical and data deception components the plant and anomaly detector are linear systems, (11) and (12) respectively, Definition 2 states that these type of attacks are 0-stealthy if $r_k^a = 0$, $k = k_0, \dots, k_f$. The idea of 0-stealthy attacks then consists of

designing an attack policy and attack signal $\mathcal{A}_{k_0}^{k_f}$ so that the residue r_k does not change due to the attack.

A particular subset of 0-stealthy attacks are characterized in the following lemma:

LEMMA 1. *The attack signal $\mathcal{A}_{k_0}^{k_f}$ is 0-stealthy with respect to any \mathcal{D} if $\tilde{y}_k^a = 0$, $\forall k \geq k_0$.*

PROOF. Consider the attacked components of the controller and the anomaly detector in (11) and (12) with $\hat{x}_0^a = \xi_{0|0}^a = 0$. From the controller dynamics it directly follows that $\tilde{y}_k^a = 0$, $\forall k \geq k_0$ results in $u_k^a = 0$, $\forall k \geq k_0$, as the input to the controller (\tilde{y}_k^a) is zero. Since $\hat{x}_0^a = 0$ and $\tilde{y}_k^a = u_k^a = 0$, $\forall k \geq k_0$, meaning that the detector's inputs are zero, we then conclude $r_k^a = 0$, $\forall k \geq k_0$. \square

Both the definition of 0-stealthy attacks and Lemma 1 indicate that these attacks are decoupled from the outputs of linear systems, r_k and y_k respectively. Hence finding 0-stealthy attack signals relates to the output zeroing problem or zero dynamics studied in the control theory literature [20]. Note that such attack requires the perfect knowledge of the plant dynamics P and the attack signal is then based on the open-loop prediction of the output changes due to the attack, as illustrated in Figure 4 where \mathcal{K}_z denote the zero dynamics and there is no disclosure of sensor or actuator data.

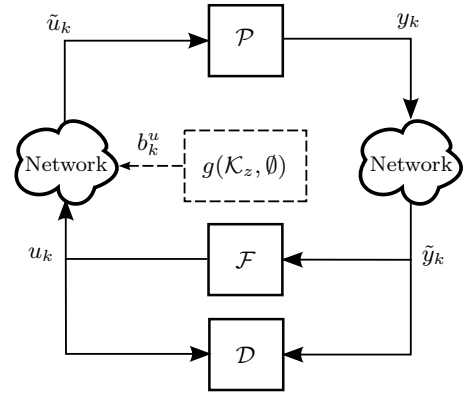


Figure 4: Schematic of the zero dynamics attack.

In Section 5 a particular instance of this attack was considered, where only the actuator data is corrupted. The zero attack policy thus corresponds to the transmission zero dynamics of the plant, which is now described. The plant dynamics due to an attack on the actuator data are described by

$$\begin{aligned} x_{k+1}^a &= Ax_k^a + Ba_k \\ \tilde{y}_k^a &= Cx_k^a \end{aligned} \quad (15)$$

with $a_k = b_k^u$. Given the discrete-time system (15) with B having full column rank, the transmission zeros can be calculated as the values $\nu \in \mathbb{C}$ that cause the matrix $P(\nu)$ to lose rank, where

$$P(\nu) = \begin{bmatrix} \nu I - A & -B \\ C & 0 \end{bmatrix}. \quad (16)$$

Those values are called minimum phase or non-minimum phase zeros depending on whether they are stable or unstable zeros, respectively. In discrete-time systems a zero is stable if $|\nu| < 1$ and unstable otherwise.

Attack policy

The attack policy then corresponds to the input sequence (a_k) that makes the outputs of the process (\tilde{y}_k^a) identically zero for all k and is illustrated in Figure 4. It can be shown [20] that the solution to this problem is given by the sequence

$$a_k = g\nu^k, \quad (17)$$

where g is the input zero direction for the chosen zero ν . The input zero direction can be obtained by solving the following equation

$$\begin{bmatrix} \nu I - A & -B \\ C & 0 \end{bmatrix} \begin{bmatrix} x_0 \\ g \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}. \quad (18)$$

where x_0 is the initial state for which the input sequence (17) results in an identically zero output, $\tilde{y}_k^a = 0 \forall k$. If the zero is stable, that is $|\nu| < 1$, the attack will asymptotically decay to zero, thus having little effect on the plant. However, in the case of unstable zeros the attack grows geometrically, which could cause a great damage to the process.

Disclosure resources

This attack scenario considers an open-loop attack policy and so no disclosure capabilities are required, resulting in $\mathcal{R}_C^u = \mathcal{R}_C^y = \emptyset$ and $\mathcal{I}_k^u = \mathcal{I}_k^y = \emptyset \forall k$.

Disruption resources

The disruption capabilities in this attack scenario correspond to the ability of performing deception attacks on the actuator data channels. Therefore the required resources are $\mathcal{R}_I^u = \{1, \dots, q\}$, $\mathcal{R}_I^y = \emptyset$, and $F = 0$

System knowledge

The ability to compute the open-loop attack policy (17) requires the perfect knowledge zero dynamics (18), which we denote as \mathcal{K}_z . Note that computing the zero dynamics requires perfect knowledge of the plant dynamics, namely A , B , and C . No knowledge of the feedback controller or anomaly detector is assumed in this scenario.

Stealthiness constraint

Note that the transmission zero attack is 0-stealthy only if $x_0^a = x_0$. However the initial condition of the system under attack x_0^a is defined to be zero at the beginning of the attack. Therefore stealthiness of the attack may be violated for large differences between $x_0^a = 0$ and x_0 .

4.4 Bias Injection Attack

Here a particular scenario of false-data injection is considered, where the attacker's goal is to inject a constant bias in the system without being detected. For this scenario, the class of α -stealthy attacks is characterized at steady-state and a method to evaluate the corresponding impact is proposed. Furthermore, we derive the policy yielding the largest impact on the system.

Denote a_∞ as the bias to be injected and recall the anomaly detector dynamics under attack given by (7). The steady-state detectability of the attack is then dependent on the steady-state value of the residual

$$r_\infty^a = (C_e(I - A_e)^{-1}B_e + D_e)a_\infty =: G_{ra}a_\infty. \quad (19)$$

The largest α -stealthy attacks are then characterized by

$$\|G_{ra}a_\infty\|_2 = \delta_\alpha. \quad (20)$$

Although attacks satisfying (20) could be detected during the transient, incipient attack signals slowly converging to a_∞ may go undetected, as it is shown in the experiments in Section 5.

The impact of such attacks can be evaluated using the closed-loop dynamics under attack given by (6). Recalling that $\eta_k = [x_k^\top \ z_k^\top]^\top$, the steady-state impact on the state is given by

$$x_\infty^a = [I \ 0](I - A_c)^{-1}B_c a_\infty =: G_{xa}a_\infty. \quad (21)$$

Largest 2-norm state bias. The α -stealthy attack yielding the largest bias in the 2-norm sense can be computed by solving

$$\max_{a_\infty} \|G_{xa}a_\infty\|_2^2 \quad (22)$$

$$\text{s.t.} \quad \|G_{ra}a_\infty\|_2^2 \leq \delta_\alpha^2.$$

Note that this problem is unbounded unless

$$\ker(G_{ra}) \subseteq \ker(G_{xa}),$$

where $\ker(A)$ denotes the null space of A , and the solution is trivial. Therefore in this section we consider the non-trivial case in which the previous condition holds.

The above optimization problem can be transformed into a generalized eigenvalue problem and the corresponding optimal solution characterized in terms of generalized eigenvalues and eigenvectors. Denote λ^* and v^* as the largest generalized eigenvalue and corresponding unit-norm eigenvector of the matrix pencil $G_{xa}^\top G_{xa} - \lambda G_{ra}^\top G_{ra}$, satisfying

$$(G_{xa}^\top G_{xa} - \lambda^* G_{ra}^\top G_{ra})v^* = 0.$$

It can be shown that the optimal solution to the optimization problem (22) is given by

$$a_\infty^* = \frac{\delta_\alpha}{\|G_{ra}v^*\|_2} v^*, \quad (23)$$

and the corresponding optimal value is $\|G_{xa}a_\infty\|_2^2 = \lambda^* \delta_\alpha^2$.

Largest infinity-norm state bias. Similarly, the α -stealthy attack yielding the largest bias in the infinity-norm sense is the solution to the following optimization problem

$$\max_{a_\infty} \|G_{xa}a_\infty\|_\infty \quad (24)$$

$$\text{s.t.} \quad \|G_{ra}a_\infty\|_2 \leq \delta_\alpha.$$

A possible method to solve this problem is to observe that

$$\|G_{xa}a_\infty\|_\infty = \max_i \|e_i^\top G_{xa}a_\infty\|_2,$$

where the vector e_i is i -th column of the identity matrix. Thus one can transform the optimization problem (24) into a set of problems with the same structure as (22), obtaining

$$\max_i \max_{a_\infty^i} \|e_i^\top G_{xa}a_\infty^i\|_\infty \quad (25)$$

$$\text{s.t.} \quad \|G_{ra}a_\infty^i\|_2 \leq \delta_\alpha.$$

Denote λ_i^* as the largest generalized eigenvalue of the matrix pencil $G_{xa}^\top e_i e_i^\top G_{xa} - \lambda G_{ra}^\top G_{ra}$. Letting $\lambda^* = \max_i \lambda_i^*$ and v^* be the corresponding generalized eigenvector, the optimal attack is given by (23) and the corresponding optimal value is $\|G_{xa}a_\infty\|_\infty = \sqrt{\lambda^*} \delta_\alpha$.

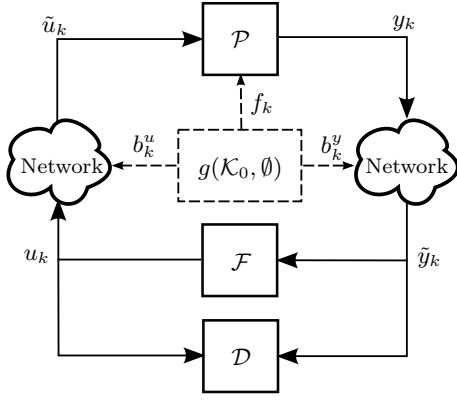


Figure 5: Schematic of the bias injection attack.

Attack policy

The bias injection attack is illustrated in Figure 5. The steady-state attack policy yielding the maximum impact on the physical system is given by (23). For the transient, we consider that the attacker uses a linear low-pass filter so that the data corruptions are slowly converging to the steady-state values. As an example, for a set of identical first-order filters the open-loop attack sequence is described by

$$a_{k+1} = \beta a_k + (1 - \beta) a_\infty^*, \quad (26)$$

where $0 < \beta < 1$ and $a_0 = 0$.

Disclosure resources

Similarly to the zero attack, no disclosure capabilities are required for this attack, since the attack policy is open-loop. Therefore we have $\mathcal{R}_C^u = \mathcal{R}_C^y = \emptyset$ and $\mathcal{I}_k^u = \mathcal{I}_k^y = \emptyset \forall k$.

Disruption resources

The biases may be added to both the actuator and sensor data, hence the required resources are $\mathcal{R}_I^u \subseteq \{1, \dots, q\}$, $\mathcal{R}_I^y \subseteq \{1, \dots, p\}$. Since no physical attack is performed, we have $F = 0$.

System knowledge

As seen in (22), the open-loop attack policy (26) requires the knowledge of the closed-loop system and anomaly detector steady-state gains G_{ra} and G_{xa} , which we denoted as \mathcal{K}_0 as shown in Figure 5.

Stealthiness constraint

Note that the steady-state value of the data corruption a_∞^* is only necessary for the attack to be α -stealthy, since the transients are disregarded. In practice, however, it has been observed in the Fault Diagnosis literature that incipient faults with slow dynamics are hard to detect [5]. Therefore the low-pass filter dynamics in (26) could be designed sufficiently slow as to difficult detection.

5. EXPERIMENTS

In this section we present our testbed and report experiments on staged cyber attacks following the different scenarios described in the previous section.

5.1 Quadruple-Tank Process

Our testbed consists of a Quadruple-Tank Process (QTP) [11] controlled through a wireless communication network, as shown in Figure 6.

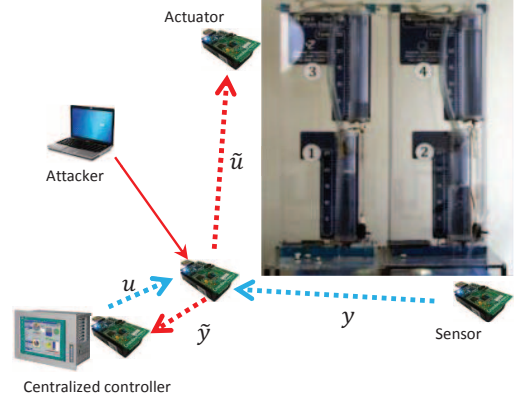


Figure 6: Schematic diagram of the testbed with the Quadruple-Tank Process and a multi-hop communication network.

The plant model can be found in [11]

$$\begin{aligned} \dot{h}_1 &= -\frac{a_1}{A_1} \sqrt{2gh_1} + \frac{a_3}{A_1} \sqrt{2gh_3} + \frac{\gamma_1 k_1}{A_1} u_1, \\ \dot{h}_2 &= -\frac{a_2}{A_2} \sqrt{2gh_2} + \frac{a_4}{A_2} \sqrt{2gh_4} + \frac{\gamma_2 k_2}{A_2} u_2, \\ \dot{h}_3 &= -\frac{a_3}{A_3} \sqrt{2gh_3} + \frac{(1 - \gamma_2) k_2}{A_3} u_2, \\ \dot{h}_4 &= -\frac{a_4}{A_4} \sqrt{2gh_4} + \frac{(1 - \gamma_1) k_1}{A_4} u_1, \end{aligned} \quad (27)$$

where h_i are the heights of water in each tank, A_i the cross-section area of the tanks, a_i the cross-section area of the outlet hole, k_i the pump constants, γ_i the flow ratios and g the gravity acceleration. The nonlinear plant model is linearized for a given operating point.

The QTP is controlled using a centralized controller running in a remote computer and a wireless network is used for the communications. A Kalman-filter-based anomaly detector is also running in the remote computer and alarms are triggered according to (4), for which we computed $\delta_r = 0.15$ and chose $\delta_\alpha = 0.25$ for illustration purposes. The communication network is multi-hop, having one additional wireless device relaying the data, as illustrated in Figure 6.

5.2 Replay Attack

In this scenario, the QTP is operating at a constant set-point and a hacker wants to steal water from the tank 4, the upper tank on the right side. The attacker has been able to hack one of the relay nodes that is between the sensor 2 (y_2) and the controller in a way that data from the real sensor can be recorded. Furthermore the attacker is able to replace the measurements sent to the controller with the recorded data. An example of this attack is presented in the Figure 7, where the attacker replays past data from y_2 while stealing water from tank 4. As we can see, the residue stays almost constant and therefore the attack is not detected, while there is a significant drop in water level in tanks 2 and 4.

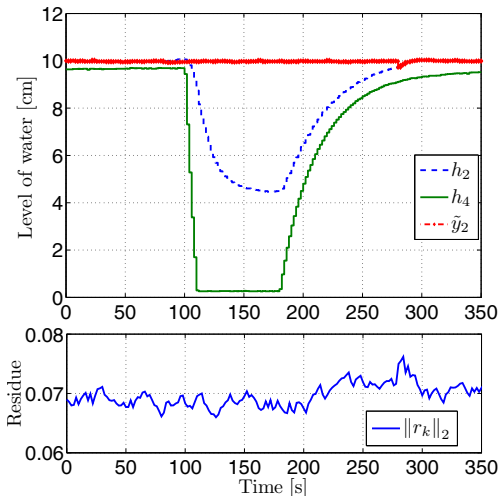


Figure 7: Results for the replay attack performed against sensor 2 from $t \approx 100s$ to $t \approx 280s$. Additionally, the attacker opens the tap of tank 4 at $t \approx 100s$ and closes it at $t \approx 180s$.

5.3 Zero Dynamics Attack

The QTP has a non-minimum phase configuration in which the plant contains an unstable zero. In this case, as discussed in Section 4.3, an attacker able to corrupt all the actuator channels may launch a false-data injection attack where the false-data follows the transmission zero dynamics, rendering the attack undetectable. This scenario is illustrated in Figure 8.

The attack remains undetected for quite some time as expected from the theory, even though the QTP is a nonlinear process. However, the fast increment of the attack signal causes saturation of the water levels after some sampling periods, as seen in Figure 8. From that moment the system dynamics change and therefore the attack signal no longer corresponds to the zero dynamics and will be detected, although it may have already damaged the system. Thus these attacks are particularly dangerous in processes that have unstable zero dynamics and in which the actuators are over-dimensioned, allowing the adversary to perform longer attacks before saturating.

5.4 Bias Injection Attack

Figure 9 shows the maximum attack impacts for all the combinations of compromised sensors and actuators in the QTP. The blue area represents the possible impacts in the process for a given number of attackable channels. As we can see, when the adversary can attack more than two channels the impact is unbounded (assuming linear dynamics), although in practice this is prevented due to saturation, as previously shown for the zero dynamics attack.

The results for the case where u_1 and y_1 are respectively corrupted with b_∞^u and b_∞^y are presented in the Figure 10, where the attacker aimed at maximizing the state bias in the infinity-norm sense while remaining stealthy. The false bias was slowly injected using a first-order low-pass filter and the following steady-state value $a_\infty = [b_\infty^u \ b_\infty^y]^\top = [2.15 \ -9.42]^\top$.

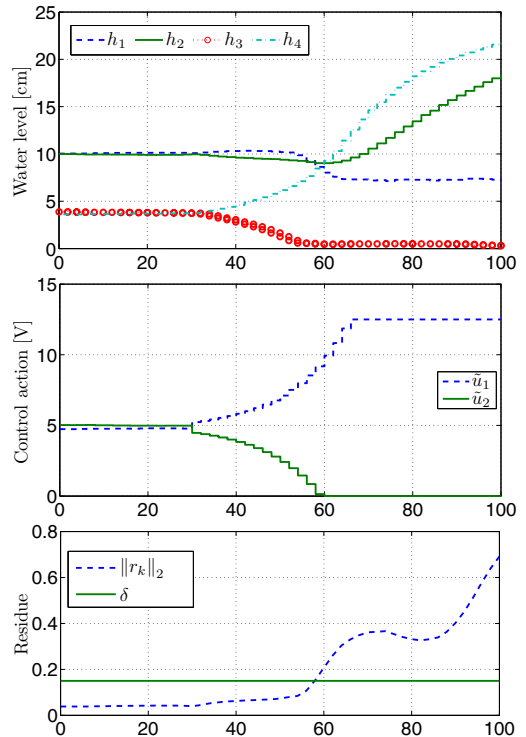


Figure 8: Results for the zero dynamics attack starting at $t \approx 30s$. Tank 3 is emptied at $t \approx 55s$, resulting in a steep increase in the residual since the linearized model is no longer valid.

6. CONCLUSIONS

In this paper we have analyzed the security of networked control systems. A novel attack space based on the attacker's system knowledge, disclosure, and disruption resources was proposed and the corresponding adversary model described. Attack scenarios corresponding to replay, zero dynamics, and bias injection attacks were analyzed using this framework. In particular the maximum impact of stealthy bias injection attacks was derived and it was shown that the corresponding policy does not require perfect model knowledge. These attack scenarios were illustrated using an experimental setup based on a quadruple-tank process controlled over a wireless network.

7. ACKNOWLEDGMENTS

This work was supported in part by the European Commission through the VIKING project, the Swedish Research Council under Grants 2007-6350 and 2009-4565, and the Knut and Alice Wallenberg Foundation.

8. REFERENCES

- [1] S. Amin, A. Cárdenas, and S. Sastry. Safe and secure networked control systems under denial-of-service attacks. In *Hybrid Systems: Computation and Control*, pages 31–45. Lecture Notes in Computer Science. Springer Berlin / Heidelberg, April 2009.
- [2] S. Amin, X. Litrico, S. S. Sastry, and A. M. Bayen. Stealthy deception attacks on water scada systems. In

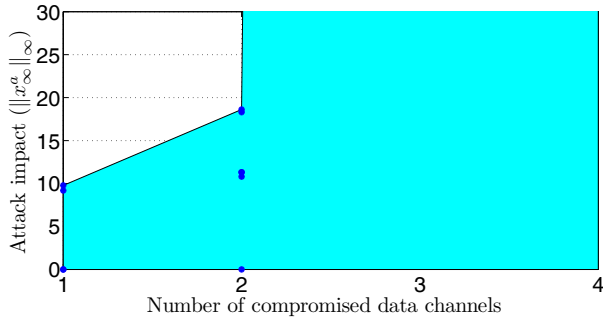


Figure 9: Maximum impact of stealthy bias injection attacks in the QTP for all sets of compromised sensors and actuators.

Proc. of the 13th ACM Int. Conf. on Hybrid systems: computation and control, HSCC '10, New York, NY, USA, 2010. ACM.

- [3] M. Bishop. *Computer Security: Art and Science*. Addison-Wesley Professional, 2002.
- [4] A. Cárdenas, S. Amin, and S. Sastry. Research challenges for the security of control systems. In *Proc. 3rd USENIX Workshop on Hot topics in security*, July 2008.
- [5] J. Chen and R. J. Patton. *Robust Model-Based Fault Diagnosis for Dynamic Systems*. Kluwer Academic Publishers, 1999.
- [6] S. X. Ding. *Model-based Fault Diagnosis Techniques: Design Schemes*. Springer Verlag, 2008.
- [7] A. Giani, S. Sastry, K. H. Johansson, and H. Sandberg. The VIKING project: an initiative on resilient control of power networks. In *Proc. 2nd Int. Symp. on Resilient Control Systems*, Idaho Falls, ID, USA, Aug. 2009.
- [8] S. Gorman. Electricity grid in U.S. penetrated by spies. *The Wall Street Journal*, page A1, April 8th 2009.
- [9] A. Gupta, C. Langbort, and T. Başar. Optimal control in the presence of an intelligent jammer with limited actions. In *Proc. of the 49th IEEE Conf. on Decision and Control (CDC)*, Dec. 2010.
- [10] I. Hwang, S. Kim, Y. Kim, and C. E. Seah. A survey of fault detection, isolation, and reconfiguration methods. *IEEE Transactions on Control Systems Technology*, 18(3):636–653, May 2010.
- [11] K. Johansson. The quadruple-tank process: a multivariable laboratory process with an adjustable zero. *IEEE Transactions on Control Systems Technology*, 8(3):456–465, May 2000.
- [12] O. Kosut, L. Jia, R. Thomas, and L. Tong. Malicious data attacks on smart grid state estimation: Attack strategies and countermeasures. In *Proc. of IEEE SmartGridComm*, Oct. 2010.
- [13] Y. Liu, M. K. Reiter, and P. Ning. False data injection attacks against state estimation in electric power grids. In *Proc. 16th ACM Conf. on Computer and Communications Security*, pages 21–32, New York, NY, USA, 2009.
- [14] Y. Mo and B. Sinopoli. Secure control against replay

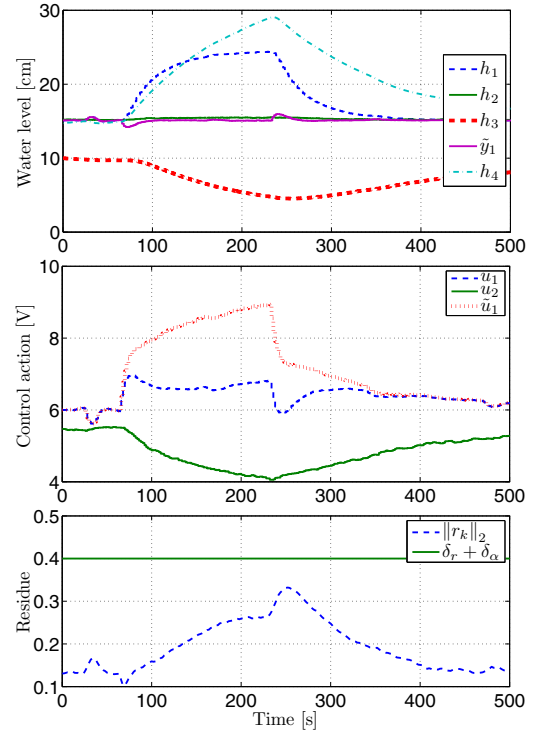


Figure 10: Results for the bias attack against the actuator 1 and sensor 1 in the minimum phase QTP. The attack is launched using a low-pass filter in the instant $t \approx 70s$ and stopped at $t \approx 230s$.

attack. In *47th Annual Allerton Conference on Communication, Control, and Computing*, Oct. 2009.

- [15] F. Pasqualetti, F. Dorfler, and F. Bullo. Cyber-physical attacks in power networks: Models, fundamental limitations and monitor design. In *Proc. of the 50th IEEE Conf. on Decision and Control and European Control Conference*, Orlando, FL, USA, Dec. 2011.
- [16] H. Sandberg, A. Teixeira, and K. H. Johansson. On security indices for state estimators in power networks. In *Preprints of the First Workshop on Secure Control Systems, CPSWEEK 2010*, Stockholm, Sweden, April 2010.
- [17] R. Smith. A decoupled feedback structure for covertly appropriating networked control systems. In *Proc. of the 18th IFAC World Congress*, Milano, Italy, August-September 2011.
- [18] A. Teixeira, H. Sandberg, G. Dán, and K. H. Johansson. Optimal power flow: Closing the loop over corrupted data. In *Proc. American Control Conference*, 2012. Accepted.
- [19] L. Xie, Y. Mo, and B. Sinopoli. False data injection attacks in electricity markets. In *First IEEE International Conference on Smart Grid Communications*, Oct. 2010.
- [20] K. Zhou, J. C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1996.