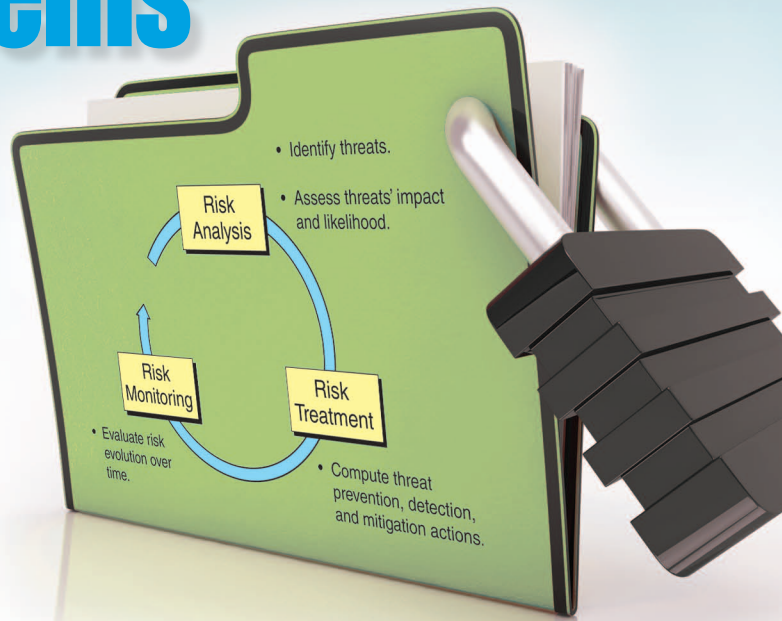


Secure Control Systems

A QUANTITATIVE RISK MANAGEMENT APPROACH

ANDRÉ TEIXEIRA,
KIN CHEONG SOU,
HENRIK SANDBERG, and
KARL HENRIK JOHANSSON



Critical infrastructures must continuously operate safely and reliably, despite a variety of potential system disturbances. Given their strict operating requirements, such systems are automated and controlled in real time by several digital controllers receiving measurements from sensors and transmitting control signals to actuators. Since these physical systems are often spatially distributed, there is a need for information technology (IT) infrastructures enabling the timely data flow between the system components. These networked control systems are ubiquitous in modern societies [1]. Examples include the electric power network, intelligent transport systems, and industrial processes.

Networked control systems are vulnerable to cyberthreats through the use of open communication networks and heterogeneous IT components. Because networked control systems are often operated through supervisory control and data acquisition (SCADA) systems and the measurement and control data in these systems are commonly transmitted through unprotected communication channels, the networked control system is vulnerable to several threats [2]. As illustrative examples, we mention the cyberattacks on power transmission networks operated by SCADA systems [3] and the Stuxnet malware that allegedly

infected an industrial control system and disrupted its operation [4], [5].

Unlike other IT systems where cybersecurity mainly involves the protection of data, cyberattacks on networked control systems may influence physical processes through feedback actuation. Therefore, networked control-system security needs to consider threats at both the cyber and physical layers.

Control theory has developed frameworks to handle disturbances and faults [6], [7], and these tools can be used to detect and attenuate the consequences of cyberattacks on networked control systems. However, there are substantial conceptual and technical differences between a fault-tolerant and a secure control framework. Faults are commonly considered to be physical events that affect the system behavior. Simultaneous events are assumed to be noncolluding, in the sense that events do not act in a coordinated way. On the other hand, cyberattacks may be performed over a significant number of attack points in a coordinated fashion; see, for instance, [8]–[10]. Moreover, faults are constrained by the physical dynamics and do not have an intent or objective to fulfill, as opposed to cyberattacks that do have a malicious intent and are not directly constrained by the dynamics of the physical process. Ensuring security may involve addressing a large number of threats, thus requiring the use of risk management methods [11] to prioritize the threats to be mitigated.

The need for novel methods to enhance the cybersecurity of networked control systems has motivated several research directions recently. The general problem of networked control systems under cyberattack is discussed and formalized in [12] and [13]. Various attack scenarios are evaluated on real and simulated benchmark systems in [13] and [14], respectively. For electric power networks, false-data injection attacks have been analyzed in detail in terms of vulnerability quantification [15], attack impact [16], [17], detection schemes [18], and attack evaluation on a realistic test bed [8]. Specific classes of attacks have been analyzed for dynamic control systems, such as replay attacks [19], stealthy false-data injection attacks [9], [10], [13], and denial-of-service attacks [20]. Quantification of tolerable errors in sensor and actuator data and their mitigation is discussed in [21], while [22]–[24] studied the network-wide data dissemination under malicious links, and [25] proposed a game-theoretic approach to cross-layer security under cascading failures.

This article presents some of the recent approaches to address cybersecurity of networked control systems under the unified perspective of risk management. First the architecture and modeling assumptions of the networked control system and adversary are introduced, following the work in [13]. Specifically, we describe the models and assumptions used for the plant, communication network, controller, and anomaly detector. Moreover, important concepts regarding the system's operation are defined, such as the system's nominal behavior and safe sets. The adversary's model, goals, and constraints are also discussed. After describing three fundamental security properties of IT systems, the adversary model is defined in terms of available resources to violate the aforementioned properties, knowledge of the system's model, and a given attack policy. The attack policy is designed according to the adversary's aims: to produce the maximum impact on the physical plant while remaining stealthy.

To tackle the existing threats, a defense methodology based on the risk management framework is presented. The notion of risk is defined in terms of a threat's scenario, impact, and likelihood, and the risk management framework is described. In this article, emphasis is given to the assessment and treatment of risk. In particular, recent quantitative tools developed in [26] and [27] for analyzing the risk of threats of static and dynamic systems are presented. The risk assessment method from [26] is tailored to quantify the likelihood of threats on a static electric power system, while the approach in [27] addresses dynamic systems and analyzes both the likelihood and impact of threats. The proposed risk assessment methods attempt to quantify the risk of different hypothetical attack scenarios for the present configuration and model of the system. As such, these methods are not executed based on real-time data of the system. The outcome from the risk assessment methods may be used for risk treatment, which is also discussed in this article and related to previous work for static and dynamic systems, [18], [28], [29] and [19], [30], respectively.

The outline of this article is as follows. First, the models for the networked control system and adversary are discussed in "Networked Control Systems Under Attacks," together with the risk management framework. The article proceeds by presenting the risk analysis methods in "Risk Analysis for Stealthy Deception Attacks." First, the method for static systems is described in detail and illustrated for large-scale electric power systems. Risk treatment methodologies for electric power systems are also discussed. Next, the risk assessment method for the dynamic case considering impact and likelihood is briefly presented and illustrated on a wireless quadruple-tank test bed. Possible risk treatment approaches are also discussed and illustrated. A summary of the article and concluding remarks are presented in the last section.

NETWORKED CONTROL SYSTEM UNDER ATTACKS

Networked control systems are spatially distributed systems where the physical plant is operated by digital controllers that receive measurements from spatially distributed sensors and transmit control signals to spatially distributed actuators through a communication network; see [1] and references therein. A typical networked control system structure has the four main components given in Figure 1: the physical plant, communication network, digital feedback controller, and digital anomaly detector. A model of each component is described below. Since this article concerns mainly threats arising from the cyber side of the networked control system, only discrete-time models are discussed in this article, where $k \in \mathbb{N}$ is the integer time index. Relevant work on the security of networked control system using continuous-time models may be found in [9] and [10].

The plant operation is supported by a communication network through which the sensor measurements and

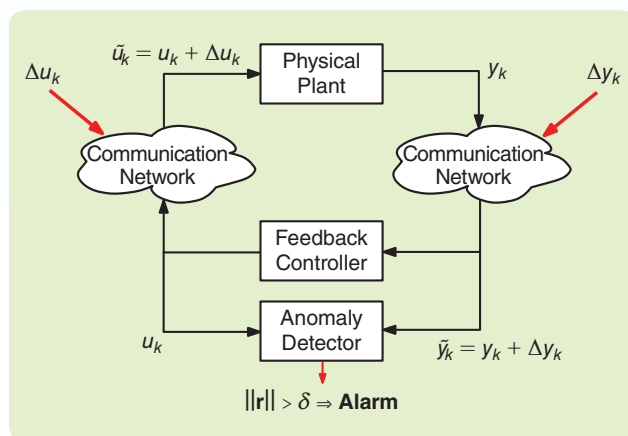


FIGURE 1 A schematic of a networked control system under attack. The plant exchanges data with the feedback controller and anomaly detector through a communication network, where \tilde{u}_k (u_k) and y_k (\hat{y}_k) are the control and measurement signals on the plant (controller) side. An adversary may inject false data Δu_k and Δy_k through the communication network. An alarm is triggered by the anomaly detector when the norm of the residue signal r over the time interval $[k_0, k_f]$ exceeds a given threshold δ .

actuator data are transmitted. On the plant side, the measurement data correspond to $y_k \in \mathbb{R}^{n_y}$, while $\tilde{u}_k \in \mathbb{R}^{n_u}$ represents the actuator data. On the controller side, the sensor and actuator data are denoted by $\tilde{y}_k \in \mathbb{R}^{n_y}$ and $u_k \in \mathbb{R}^{n_u}$, respectively. The dynamical model of the plant is given by

$$\begin{cases} x_{k+1} = f_x(x_k, \tilde{u}_k, d_k), \\ y_k = g_x(x_k, \tilde{u}_k, d_k), \end{cases} \quad (1)$$

where $x_k \in \mathbb{R}^{n_x}$ denotes the plant's state, and the unknown input $d_k \in \mathbb{R}^{n_d}$ models possible disturbances or faults affecting the system. Mismatches between the transmitted signals, y_k and u_k , and the received ones, \tilde{y}_k and \tilde{u}_k , may be due to the communication network, as in the case of delays or packet drops. Since the focus of this article is on the security of the control system with respect to malicious adversaries, the communication network is assumed to be ideal and process and measurement noises are neglected. Under this assumption, the mismatches between the transmitted signals and the received ones are caused by the adversary's actions. Similarly, physical attacks performed by malicious adversaries are modeled by the unknown input d_k . The nominal behavior of the system under no attack is defined as follows.

Definition 1

A networked control system is said to have *nominal behavior* if $\tilde{u}_k = u_k$, $\tilde{y}_k = y_k$, and $d_k = 0$. Otherwise, the system has *abnormal behavior*.

Several physical systems have tight operating constraints that, if not satisfied, might result in physical damage to the system, for example, power systems, where electrical power flows along transmission lines cannot exceed physical limits. In this work, the concept of safe sets is used to characterize the safety constraints. Consider the time interval $[0, N]$ and define the vector $\mathbf{x} \triangleq [x_0^\top \dots x_N^\top]^\top \in \mathbb{R}^{n_x(N+1)}$. Given the set $\mathcal{S}_x \subseteq \mathbb{R}^{n_x(N+1)}$, safety is defined as follows.

Definition 2

The system is said to be *safe* over the time interval $[0, N]$ if the state trajectory \mathbf{x} is contained in the *safe set* \mathcal{S}_x .

Returning to the power system example, let x_k be the state of the power system and denote the output $y_k = C_x x_k$ as the instantaneous power flow measured on a given transmission line. Due to physical limits, the cable cannot sustain an arbitrarily large instantaneous power. With the appropriate scaling of y_k , such an operating constraint can be defined in terms of the safe set $\mathcal{S}_x = \{\mathbf{x}: \max_k \|C_x x_k\|_\infty \leq 1\}$.

To comply with performance requirements in the presence of unknown disturbances, the physical plant is assumed to be controlled by an appropriate feedback controller [6], which computes the control signal u_k given the measurement signal \tilde{y}_k received through the communication network. The output feedback controller can be written in a state-space form as

$$\begin{cases} z_{k+1} = f_z(z_k, \tilde{y}_k), \\ u_k = g_z(z_k, \tilde{y}_k), \end{cases} \quad (2)$$

where the states of the controller are labeled as $z_k \in \mathbb{R}^{n_z}$. Given the plant model, the controller is supposed to be designed so that acceptable performance is achieved under nominal behavior.

The anomaly detector, which is collocated with the controller, monitors the system to detect possible deviations from the nominal behavior. It has access to only \tilde{y}_k and u_k . Several approaches to detecting malfunctions in control systems are available in the fault diagnosis literature [7], [31]. Other schemes tailored to detecting sparse adversarial attacks have also been proposed [32], [33]. A common approach is the observer-based fault detection filter

$$\begin{cases} s_{k+1} = f_s(s_k, u_k, \tilde{y}_{k+1}), \\ r_k = g_s(s_k, u_k, \tilde{y}_{k+1}), \end{cases} \quad (3)$$

where $s_k \in \mathbb{R}^{n_s}$ is the anomaly detector's state. Based on the plant and controller models, the control signal u_k , and the received measurements \tilde{y}_k , the fault detection filter computes the residue $r_k \in \mathbb{R}^{n_r}$. The residue signal is evaluated to detect and locate existing anomalies, as depicted in Figure 1.

The anomaly detector (3) is designed such that

- 1) under nominal behavior of the system ($u_k = \tilde{u}_k$, $y_k = \tilde{y}_k$), the expected value of r_k converges asymptotically to a neighborhood of the origin
- 2) the residue is sensitive to the anomalies so that an abnormal behavior of the system results in a nonzero residue signal.

The aim of an anomaly detector is to evaluate the residue to detect anomalies with high probability while keeping the rate of false alarms due to uncertainties below a certain level. Given the aforementioned design specifications of the anomaly detector, various residue evaluation techniques described in [34] may be used to detect anomalies. For instance, the anomaly detector may be designed to trigger an alarm when a given norm of the residue signal exceeds a certain bound over the time interval $[k_0, k_f]$

$$\|\mathbf{r}\|_p \triangleq \left(\sum_{i=1}^{(k_f - k_0 + 1)n_r} |r_i|^p \right)^{\frac{1}{p}} > \delta, \quad (4)$$

where $\mathbf{r} = [r_{k_0}^\top \dots r_{k_f}^\top]^\top$ is the residue signal over the time interval $[k_0, k_f]$, $\|\mathbf{r}\|_p$ denotes the ℓ_p -norm of \mathbf{r} , and $\delta > 0$ ensures a desired false-alarm rate with respect to uncertainties. Different methods to compute the threshold δ and the corresponding false-alarm rate may be found in [31].

Adversary Model

Due to the tight coupling between the cyber and physical domains, the control system behavior depends on the state and properties of the IT infrastructure. To model and understand how a cyber adversary may affect the networked control system operation requires knowing how IT systems are vulnerable to adversaries. Computer security literature identifies three fundamental properties of information and services in IT systems, namely *confidentiality*, *integrity*, and *availability*, often denoted as CIA [35]. They

can be violated by disclosure, deception, and denial-of-service attacks, respectively. For examples of attacks violating these properties in networked control systems, see “The CIA in Networked Control Systems.”

Disclosure attacks enable the adversary to gather sequences of data \mathcal{I}_k from the calculated control actions u_k and the real measurements y_k . As such, the physical dynamics of the system are not affected by this type of attack. Instead, these attacks gather intelligence that may enable more complex attacks, such as replay attacks [36]. On the other hand, 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. The deception attacks are modeled as

$$\begin{aligned}\tilde{u}_k &\triangleq u_k + \Delta u_k, \\ \tilde{y}_k &\triangleq y_k + \Delta y_k,\end{aligned}$$

where the vectors Δu_k and Δy_k represent the data corruption to the respective data channels, as depicted in Figure 1. The data corruption vectors may have sparsity patterns according to the adversary’s resources, namely the communication channels that can be corrupted. Similarly, denial-of-service attacks may also affect the transmitted data by preventing it from reaching the desired destination. Attacks that may affect the system behavior directly and through feedback are classified as disruption attacks [13]. From the preceding discussion, we conclude that physical, deception, and denial-of-service attacks are classified as disruption attacks. The data channels and physical actuators required to perform specific disclosure and disruption attacks are denoted as disclosure and disruption resources, respectively.

In addition to the disclosure and disruption resources required to stage a given attack, the adversary’s resources can also include knowledge of the system model. Different attack scenarios can be qualitatively categorized in terms of the required resources in the attack space, as illustrated in Figure 2. A given point in the attack space represents an instance of the adversary model in Figure 3 where each of the adversary resources is mapped to a specific axis of the attack space. The attack policy mapping the model knowledge \mathcal{K} and disclosed data gathered until time k , \mathcal{I}_k , to the attack vector $a_k \in \mathbb{R}^n$ is denoted as $a_k = g(\mathcal{K}, \mathcal{I}_k)$.

For each attack scenario, the attack policy is designed according to the adversary’s intent, namely the attack goals and constraints. In particular, the attack scenarios in this article consider adversaries whose goal is to drive the state trajectory x of the physical system to an unsafe set while remaining stealthy, as illustrated in Figure 4. Therefore the attack goals are stated in terms of the attack impact on the system operation, while the constraints are related to the attack detectability.

The physical impact of an attack can be evaluated by assessing whether or not the state of the system remained in the safe set during and after the attack. The attack is considered successful if the state is driven out of the safe set. Attack constraints imply that attacks are constrained to remain stealthy. Denoting

$\mathbf{a} = [a_{k_0}^\top \cdots a_{k_f}^\top]^\top$ as the attack signal in the time interval $[k_0, k_f]$, and recalling that the residue signal is a function of the attack signal, a stealthy attack is defined as follows.

Definition 3

The attack signal \mathbf{a} is *stealthy* over the time interval $[k_0, k_f]$ if the magnitude of the residue signal is smaller than the detection threshold, so that no alarm is triggered.

Below, it is assumed that the disruptive attack component consists of only data deception attacks and thus at time k the attack vector $a_k = [\Delta u_k^\top \Delta y_k^\top]^\top$.

Defense Methodology

This subsection describes a common methodology to enhance a system’s cybersecurity, namely the risk management framework [35], [37], [38]. The main objective of risk

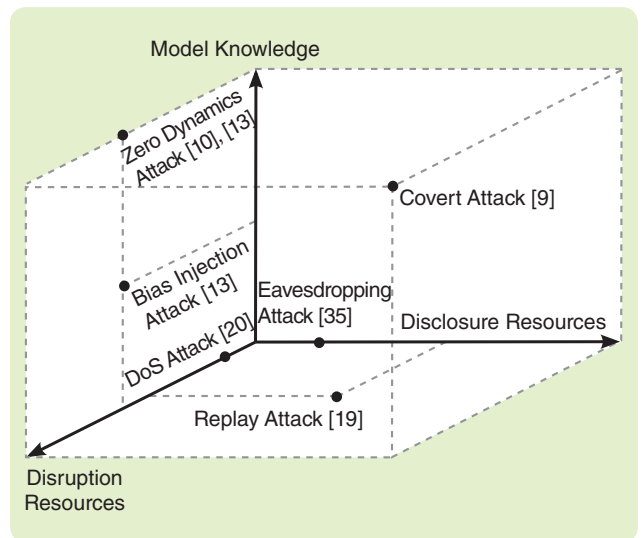


FIGURE 2 The cyberphysical attack space. Each axis of the attack space corresponds to a class of adversary resources. Several attack scenarios analyzed in related work are depicted and qualitatively categorized in the attack space.

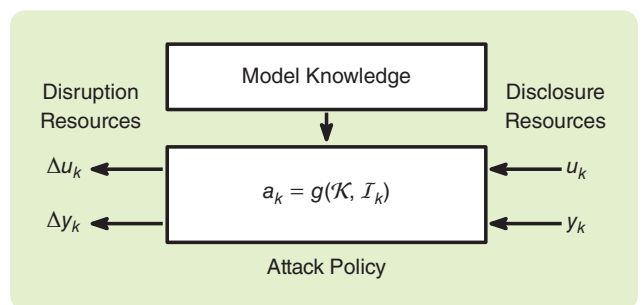


FIGURE 3 A diagram of the adversary model. The a priori model knowledge possessed by the adversary is denoted as \mathcal{K} , while \mathcal{I}_k corresponds to the set of sensor and actuator data available to the adversary, obtained through the disclosure resources, and $a_k = [\Delta u_k^\top \Delta y_k^\top]^\top$ is the attack vector that may affect the system behavior using the disruption resources. The attack policy $g(\cdot)$ maps the model knowledge and disclosed data to the attack vector.

management is to assess and minimize the risk of threats, where the notion of risk is defined as follows [39].

Definition 4

Consider a given attack *threat scenario*, the corresponding *impact* to the system, and the *likelihood* of such scenario. The

risk of the system is denoted as the set of triplets $Risk \equiv \{(Scenario, Impact, Likelihood)\}$.

The risk of different threat scenarios may be summarized in a two-dimensional risk matrix [38], where each dimension corresponds to the likelihood and impact of threats, respectively. Additionally, the risk of different threats may be compared

The CIA in Networked Control Systems

Three fundamental properties of information and services in IT systems are mentioned in the computer security literature [35] using the acronym CIA: confidentiality, integrity, and availability. Confidentiality concerns the concealment of data, ensuring it remains known only to the authorized parties. Integrity relates to the trustworthiness of data, meaning there is no unauthorized change to the information between the source and destination. Availability considers the timely access to information or system functionalities.

Figure S1 illustrates cyberattacks that violate each security property. In all three cases, the plant is sending the measurement vector $y_k = [2, 13]^T$ to the controller through a communication network. This is a private message, hence only the plant and the controller should know the message contents.

In Figure S1(a), the adversary is able to eavesdrop on the communication, thus getting access to the contents of the message. Therefore the confidentiality attribute was violated. Another scenario occurs in Figure S1(b), where the adversary succeeds in sending the false measurement vector $\tilde{y}_k = y_k + \Delta y_k$ to the controller, as if it was the plant sending it. Here data integrity is violated. In the final example, illustrated in Figure S1(c), the message sent by the plant is actually blocked and does not reach the controller. In this instance, data availability was compromised.

Whereas in IT systems the impact of such cyberattacks remains in the cyber realm, in networked control systems the impact may carry dire consequences for the physical side. Next, a specific example illustrating an attack scenario is presented, and its consequence to the physical system is discussed.

Consider a remotely controlled power generator, with θ and ω denoting its phase-angle and frequency deviation, respectively. Considering the single-machine infinite-bus model [68], the generator dynamics are described in continuous time by the normalized swing equation

$$\begin{aligned} \omega(t) &= \dot{\theta}(t), \\ M\dot{\omega}(t) &= -D\omega(t) - P_f(t) + u(t), \end{aligned}$$

where $u(t)$ is the normalized mechanical power provided to the generator and M and D are the inertia and damping coefficients, respectively. The term $P_f(t) = b \sin(\theta(t))$ corresponds to the electric power flow from the generator to the bus, where b is the susceptance parameter of the transmission line. Linearizing the model at $\omega = \theta = 0$ with $M = D = b = 1$, with the

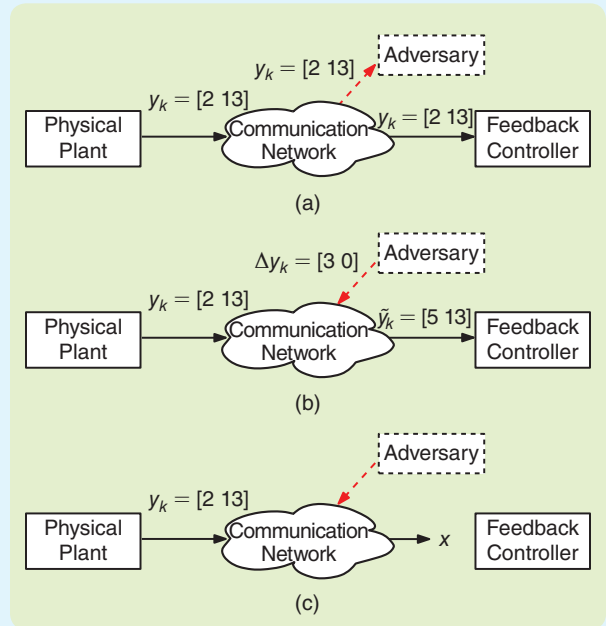


FIGURE S1 Cyberattacks on a communication network: (a) data confidentiality violation by a disclosure attack, (b) data integrity violation by a false-data injection attack, and (c) data availability violation by a denial-of-service attack.

sampling period $T_s = 1$ s, and defining the discrete-time state $x_k = [\theta_k \ \omega_k]^T$, the discrete-time model is

$$\begin{aligned} x_{k+1} &= \begin{bmatrix} 0.66 & 0.53 \\ -0.53 & 0.13 \end{bmatrix} x_k + \begin{bmatrix} 0.34 \\ 0.53 \end{bmatrix} \tilde{u}_k, \\ y_k &= C_x x_k, \end{aligned}$$

where $C_x = I$, and \tilde{u}_k is the control signal received on the plant side. Additionally, the system is safe if the frequency deviation ω is small and the power flow P_f does not exceed the line ratings. In particular, the system is said to be safe if $|\omega_k| \leq 0.05$ and $|P_k| = |\theta_k| \leq 0.1$ for all k . Defining the diagonal matrix $T = \text{diag}(10, 20)$ and $\mathbf{x} = [x_0^T \dots x_M^T]^T$ as the state trajectory over the time interval $[0, M]$, the corresponding safe set is given by

$$S_x = \{\mathbf{x}: \|T\mathbf{x}_k\|_\infty \leq 1, \forall k \in [0, M]\}.$$

The anomaly detector corresponds to the state observer

$$\begin{aligned} z_{k+1} &= (A_x - LC_x)z_k + B_x u_k + L\tilde{y}_k, \\ r_k &= \tilde{y}_k - C_x z_k, \end{aligned}$$

through increasing functions of the threat's impact and likelihood. As an example, Figure 5 illustrates a medium- and a high-risk threat with similar impacts but different likelihoods.

The risk management cycle, depicted in Figure 6, is composed of risk analysis, risk treatment, and risk monitoring. Risk analysis identifies threats and assesses the respective

likelihood and impact on the system. Threats may be identified based on historical and/or empirical data of cyberattacks and known vulnerabilities in the system [38]. The likelihood of a given threat depends on the components compromised by the adversary in a given attack scenario and their respective vulnerability. Quantitative methods

where \tilde{y}_k is the measurement received at the anomaly detector side, z_k is the estimate of x_k , the residue r_k is the output estimation error, and

$$L = \begin{bmatrix} 0.36 & 0.27 \\ -0.31 & 0.08 \end{bmatrix}$$

is the observer gain matrix design such that $A_x - LC_x$ is stable. Denoting $\mathbf{r} = [r_0^T \dots r_N^T]^T$, an alarm is triggered by the anomaly detector when $\|\mathbf{r}\|_\infty > \delta = 0.01$. The output feedback controller is given by

$$\begin{aligned} z_{k+1} &= (A_x - B_x K - LC_x)z_k + L\tilde{y}_k, \\ u_k &= -Kz_k, \end{aligned}$$

where $K = [0.0556 \ 0.3306]$ is the controller gain matrix ensuring that $A_x - B_x K$ is stable.

Consider an attack scenario where the adversary knows the exact model of the plant and is able to compromise the integrity of the control signal u_k , that is, the mechanical power supplied to the generator, and the power flow measurement $y_{1k} = P_k = \theta_k$. Defining $\tilde{u}_k = u_k + \Delta u_k$, $\tilde{y}_k = y_k + \Delta y_k$, and the attack vector $a_k = [\Delta u_k \ \Delta y_{1k}]^T$ the plant under attack is described by

$$\begin{aligned} x_{k+1} &= A_x x_k + B_x u_k + B_a a_k, \\ \tilde{y}_k &= C_x x_k + D_a a_k, \end{aligned} \quad (\text{S1})$$

with

$$B_a = \begin{bmatrix} 0.34 & 0 \\ 0.53 & 0 \end{bmatrix}, \quad D_a = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}.$$

The adversary attempts to drive the system to an unsafe state while remaining stealthy. To that end, the adversary injects an increasing piece-wise constant signal into the control input, making the generator produce more power, and thus increasing the power flow P_i along the transmission line. At the same time, the power flow increase is hidden from the controller and anomaly detector by tampering with the power flow measurement. More specifically, the attack policy is chosen as sequential instances of the zero-dynamics policy [13] during N time instants, where the attack vector is constructed as

$$\begin{aligned} \kappa_k &= \left\lfloor \frac{k+1}{N} \right\rfloor, \\ a_k &= \kappa_k \lambda^k g, \end{aligned} \quad (\text{S2})$$

with $\lambda \in \mathbb{C}$ and $g \neq 0$ being the invariant zero and the corresponding input-direction satisfying

$$\begin{bmatrix} \lambda I - A_x & -B_a \\ C_x & D_a \end{bmatrix} \begin{bmatrix} x_z \\ g \end{bmatrix} = 0.$$

In the considered attack scenario, system (S1) has the zero $\lambda = 1$ on the unit circle, input direction $g = \epsilon[-11]^T$,

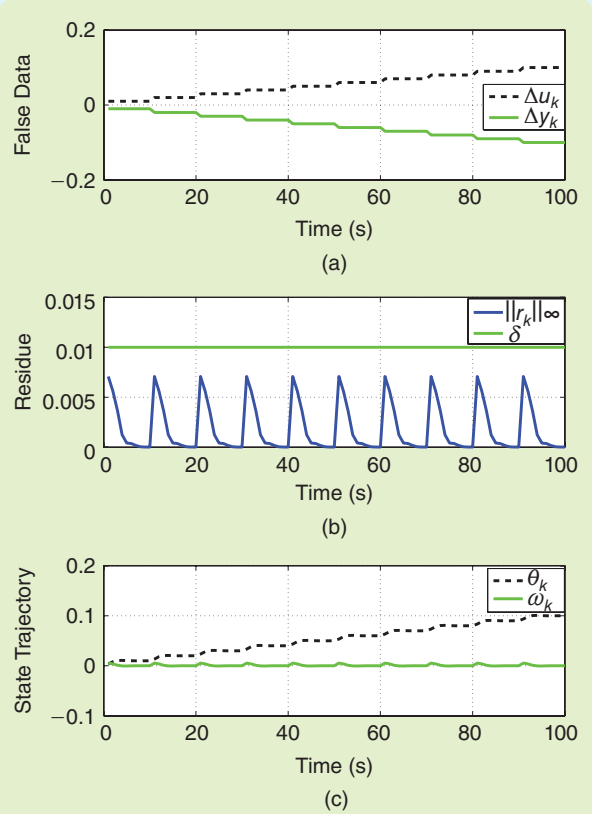


FIGURE S2 Simulation results from the attack policy (S2) with $N = 10$, $\lambda = 1$, and $g = [-0.01 \ 0.01]^T$. Plot (c) depicts the state trajectory under attack, the injected false-data is shown in (a), and (b) illustrates the corresponding residue signal and detection threshold.

and $x_z = \epsilon[-20 \ 0]^T$ for $\epsilon \neq 0$. Notice that an input of the form $a_k = \lambda^k g$ is blocked from the output at steady state by the zero $\lambda = 1$, yielding $\lim_{k \rightarrow \infty} \tilde{y}_k = 0$.

For the generator's closed-loop system with $x_0 = z_0 = 0$, choosing $N = 10$ and $\epsilon = 0.01$ and applying the attack policy (S2) results in the signals depicted in Figure S2. Observe that, at the end of the first instance, $r_N = 0$, $\theta_N = 0.01$, and $\omega_N = 0$. Therefore, the second attack instance during the time interval $[N + 1, 2N]$ begins with $r_N = 0$ and also yields $\|r_k\|_\infty \leq 0.0071$ for $k \in [N + 1, 2N]$, as shown in Figure S2. Furthermore, the final value of the power flow is increased to $P_{1N} = \theta_{2N} = 0.02$. In fact, the attack policy (S2) yields $\|r_k\|_\infty \leq 0.0071$ for all k and $P_{i,N} = \theta_{kN} = 0.01k$ for $k = 1, 2, \dots$, as depicted in Figure S2. Thus, the adversary is able to drive the system to outside the safe set at $k = 100$ while remaining stealthy.

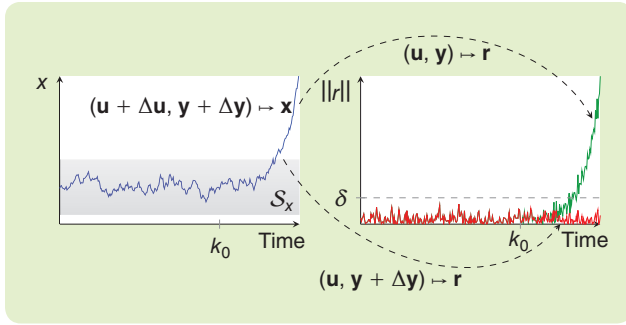


FIGURE 4 An example of state and residue signals of a networked control system under a stealthy deception attack starting at time k_0 . The plot to the left depicts the plant's state trajectory \mathbf{x} under the attacked control and measurement signals $(\mathbf{u} + \Delta\mathbf{u}, \mathbf{y} + \Delta\mathbf{y})$. The safe set S_x is indicated by the shaded region. The plot to the right depicts the instantaneous norm ($\|r\|$) of two residue signals, namely the actual residue signal (red) and the ideal one (green). The actual residue signal is computed by the anomaly detector based on the available signals $(\mathbf{u}, \mathbf{y} + \Delta\mathbf{y})$; see Figure 1. In this case, the residue norm is always smaller than δ , thus the attack is not detected while the adversary succeeds in driving the plant state out of the safe set as intended. On the other hand, if the true measurement signal \mathbf{y} is available to the anomaly detector, the residue computed from (\mathbf{u}, \mathbf{y}) successfully detects the attack.



FIGURE 5 A risk matrix plot. The threat's likelihood and impact correspond to the x axis and y axis, respectively. Two threats with a similar impact but different likelihoods are depicted. Threats with high impact and high likelihood yield a higher risk.

can be used to identify the minimal set of components that need to be compromised for each attack scenario [15], [40], while the vulnerability of each compromised components is obtained by qualitative means such as expert knowledge and historical and empirical data [40]. The potential impact of a threat may be assessed by qualitative and quantitative methods, for instance by modeling the system and simulating the attack scenarios [11].

Actions minimizing the risk of threats are determined within the risk treatment step. The different actions can be classified as prevention, detection, and mitigation. Prevention aims at decreasing the likelihood of attacks by reducing the vulnerability of the system components, for instance by encrypting the communication channels, using firewalls, and intelligent routing algorithms [28]. On the other hand, detection is an approach in which the system is continuously monitored for anomalies caused by adversary actions. Examples of detection schemes include antivirus software, network traffic analysis [41], and fault detection algorithms [31]. Once an anomaly or attack is detected, mitigation actions may be taken to disrupt and neutralize the attack, thus reducing its impact. The attack may be neutralized by replacing the compromised components or using redundant components.

The effectiveness of the defensive actions and the evolution of risk over time is evaluated throughout the risk monitoring stage. Risk monitoring continuously assesses the known and newly discovered vulnerabilities of the system, as well as the deployment of the threat mitigation actions.

Given the importance of risk analysis and risk treatment in the risk management process, the next sections illustrate and describe in detail methods that can be used for risk analysis and risk treatment in networked control systems.

RISK ANALYSIS FOR STEALTHY DECEPTION ATTACKS

Quantitative approaches to risk analysis of stealthy deception attacks are discussed in the remainder of this article. First, a simplified static case is analyzed in detail and illustrated with a power systems example. Then, the general dynamic case is presented and illustrated on a wireless quadruple-tank test bed.

Recall that the adversary aims to drive the system to an unsafe state while remaining stealthy. Additionally, the adversary also has resource constraints, in the sense that only a small number of attack points to the system are available. This section describes a framework for performing risk analysis of data deception attacks on networked control systems, where an attack is deemed less likely the more resources it requires. Particularly, the plant (1), feedback controller (2), and anomaly detector (3) are considered to be linear time-invariant systems. Defining $\eta_k = [x_k^T z_k^T s_k^T]^T$ and $a_k = [\Delta u_k^T \Delta y_k^T]^T$, the closed-loop dynamics of the networked control system driven by deception attacks are [13]

$$\begin{aligned} \eta_{k+1} &= A\eta_k + Ba_k, \\ \tilde{y}_k &= C_y\eta_k + D_y a_k, \\ r_k &= C_r\eta_k + D_r a_k. \end{aligned} \quad (5)$$

Risk Analysis for Static Models

The risk assessment in this subsection focuses on analyzing the threat's likelihood, indicated by the minimum number of sensors that need to be compromised by the adversary for a given attack scenario. The minimum number of compromised sensors is a relevant indicator of

the threat's likelihood because the sensors are often geographically distributed in networked control systems. As a result, coordinated attacks compromising multiple sensors need to be carried out simultaneously in different locations and are therefore difficult to implement.

The model in Figure 1 is simplified in two regards. First, the plant is in steady state. That is, in (5) the state vector η_k is constant for all k , so the subscript k is omitted. The second simplification is that there is no feedback control. The simplifications are made because they can lead to a more streamlined illustration of the main concept of risk assessment. In addition, the simplified structure is relevant in its own right in analyzing the cyberphysical security of power systems. The risk assessment for general dynamic models will be deferred to a later section.

The model for risk assessment is the relationship between the static plant states x and the measurements \tilde{y} received by the anomaly detector. This is described by the expression

$$\tilde{y} = C_y x + \Delta y,$$

where C_y is the measurement matrix, and Δy is the measurement data attack. In a typical static state estimation problem, such as the power network case, there are more measurements than states, and hence C_y is assumed to have full column rank [42], [43]. Based on the risk assessment model, the least-squares estimate of the states is $(C_y^T C_y)^{-1} C_y^T \tilde{y}$, and the estimate of measurements can be expressed as $C_y (C_y^T C_y)^{-1} C_y^T \tilde{y}$. Thus, the anomaly detector, which is based on measurement residue, can be described by

$$r \triangleq S \tilde{y} = (I - C_y (C_y^T C_y)^{-1} C_y^T) \tilde{y}. \quad (6)$$

Such an anomaly detector is, in general, sufficient to detect Δy in the form of a single error involving only one faulty measurement [42], [43]. However, in the face of a coordinated malicious attack on multiple measurements, the anomaly detector can fail. In particular, in [44] it was reported that an attack of the form

$$\Delta y = C_y \Delta x, \quad (7)$$

for an arbitrary Δx would not result in any additional residue in (6), apart from the residue caused by other factors, such as measurement noise. In fact, the set of stealthy deception attacks with respect to the anomaly detector (6) and a zero detection threshold is characterized by (7), and these attacks were also experimentally verified in a realistic test bed [8]. Although stealthy attacks may be obtained from (7), distinct choices of Δx may yield attack vectors Δy requiring significantly different amounts of adversary resources, in terms of the number of nonzero entries of the attack vector Δy . This number is also an indicator of the likelihood of the success of stealthy attack, as discussed earlier in this subsection. The rest of this subsection focuses

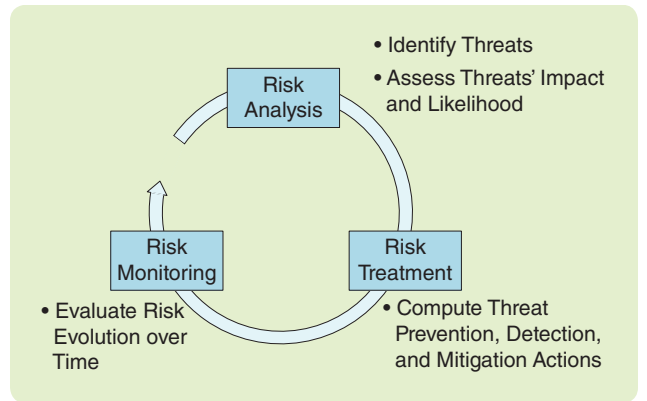


FIGURE 6 A diagram of the risk management cycle. Risk of threats is continuously minimized by iteratively performing risk analysis, risk treatment, and risk monitoring.

on the characterization of the stealthy attack vectors with the minimum number of nonzero entries, as a concrete example of the quantitative method for risk assessment.

Minimum-Resource Attacks

There is a significant amount of literature studying the stealthy attack (7) and its consequences to state-estimation data integrity (for example, [15] and [44]–[49]). It was shown numerically that stealthy attacks $\Delta y = C_y \Delta x$ are often sparse [44]. To analyze the stealthy attacks with the minimum number of nonzero entries, in [15] the notion of security index α_j for a measurement j was introduced as the optimal objective value of the following cardinality minimization problem

$$\alpha_j \triangleq \min_{\Delta x \in \mathbb{R}^n} \|C_y \Delta x\|_0 \quad \text{subject to } C_y(j,:) \Delta x \neq 0, \quad (8)$$

where $\|C_y \Delta x\|_0$ denotes the cardinality (that is, the number of nonzero entries) of the vector $C_y \Delta x$, j is the label of the measurement for which the security index α_j is computed, and $C_y(j,:)$ denotes the j th row of C_y . The security index α_j is the minimum number of measurements an attacker needs to compromise to attack measurement j without being detected by the anomaly detector. In particular, a small α_j implies that measurement j is relatively easy to compromise in a stealthy attack, therefore indicating the higher likelihood of such a threat. As a result, the knowledge of the security indices for all measurements allows the network operator to pinpoint the security vulnerabilities of the network and to better protect the network with limited resources. For example, [45] proposed a method to optimally assign limited encryption protection resources to improve the security of the network based on its security indices.

The security index (8) is a quantitative tool for risk assessment that can provide a security assessment the standard detection procedure [42], [43] might not be able to

For each attack scenario, the attack policy is designed according to the adversary's intent, namely the attack goals and constraints.

provide. As a concrete example [15], consider the measurement matrix

$$C_y = \begin{pmatrix} -1 & -1 & 0 \\ -1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & -1 \\ 0 & -1 & 0 \end{pmatrix}. \quad (9)$$

The ‘‘hat matrix’’ [42], [43], denoted K , captures how the received measurements \tilde{y} are weighted together to form a measurement estimate \hat{y} and is defined according to

$$\hat{y} = C_y \hat{x} = C_y (C_y^T C_y)^{-1} C_y^T \tilde{y} \triangleq K \tilde{y}.$$

Corresponding to the C_y in (9), is the hat matrix

$$K = \begin{pmatrix} 0.6 & 0.2 & -0.2 & 0 & 0.4 \\ 0.2 & 0.4 & -0.4 & 0 & -0.2 \\ -0.2 & -0.4 & 0.4 & 0 & 0.2 \\ 0 & 0 & 0 & 1 & 0 \\ 0.4 & -0.2 & 0.2 & 0 & 0.6 \end{pmatrix}. \quad (10)$$

The rows of the hat matrix can be used to study the measurement redundancy [42], [43]. Typically a large degree of redundancy (many nonzero entries in each row) is desirable to compensate for noisy or missing measurements. In (10), all measurements are redundant in this example except the fourth. Such nonredundant measurement is called a critical measurement. Without the critical measurement, observability is lost, meaning that it becomes impossible to uniquely determine the states based on the available measurement information. The hat matrix indicates that the critical measurement is sensitive to attacks. This is indeed the case, but some other measurements can also be vulnerable to attacks. The security indices α_j , $j = 1, \dots, 5$, respectively, are 2, 3, 3, 1, 2. Therefore, the fourth (critical) measurement has security index one, indicating that it is vulnerable to stealthy attacks. However, the first and the last measurements also have relatively small security indices. This is not obvious from K in (10). Hence, the information of the security indices can enhance the vulnerability analysis compared to the hat matrix.

Because of the cardinality minimization, computing the security indices can sometimes be hard. In fact, it can be established that problem (8) is NP-hard using techniques from [50] and [51]. As a result, known exact solution algorithms for (8) are enumerative by nature. Three different typical exact algorithms include a) enumeration on the support of $C_y \Delta x$, b) finding the maximum feasible subsystem

for an appropriately constructed system of infeasible inequalities [52], and c) the big M method (for example, [53]). This article focuses on the big M method because the resulting optimization problem can be modeled and solved using available software such as CPLEX. The big M method sets up and solves the following optimization problem.

$$\begin{aligned} & \underset{\Delta x, w}{\text{minimize}} && \sum_i w(i) \\ & \text{subject to} && C_y \Delta x \leq Mw, \\ & && -C_y \Delta x \leq Mw, \\ & && C_y(j,:) \Delta x = 1, \\ & && w(i) \in \{0, 1\} \text{ for all } i. \end{aligned} \quad (11)$$

In (11), the inequalities are interpreted entry-wise, and $0 < M < \infty$ is a user-defined constant scalar. If M is greater than the maximum entry of $C_y \Delta x^*$ in absolute value, for some optimal solution Δx^* of (8), then the optimal solution to (11) is exactly an optimal solution to (8). Otherwise, solving (11) yields a suboptimal solution, optimal among all solutions Δx such that the maximum entry of $C_y \Delta x$ is less than or equal to M in absolute value. The procedure described in [54] can always find a sufficiently large M to ensure that the big M method indeed provides the optimal solution to (8). In addition, the physics and insights of the underlying application problem can also lead to a suitable M . The optimization problem in (11) is a mixed integer linear programming (MILP) problem; see ‘‘Mixed Integer Linear Programming’’ for additional details.

For large-scale system analysis, it might be deemed impractical to obtain the exact solution to the security index problem in (8). In this case, it might be necessary to settle for an approximate solution instead. A particular method to obtain an approximate solution is ℓ_1 relaxation. For general information about ℓ_1 relaxation; see, for example, [55]–[57]. Here the properties that are most relevant to this article are described. Instead of solving (8), the ℓ_1 relaxation method sets up and solves the following optimization problem:

$$\begin{aligned} & \underset{\Delta x \in \mathbb{R}^n}{\text{minimize}} && \|C_y \Delta x\|_1 \\ & \text{subject to} && C_y(j,:) \Delta x = 1, \end{aligned} \quad (12)$$

where $\|C_y \Delta x\|_1$ denotes the vector ℓ_1 -norm (sum of absolute values of the entries) of $C_y \Delta x$. In addition, the right-hand side of the constraint in (12) needs to be normalized to ensure that the problem is well posed. Problem (12) can be

Mixed Integer Linear Programming

A mixed integer linear programming (MILP) problem is an optimization problem over both real and integer decision variables with a linear objective function and linear constraints. It is basically an LP problem except that some of the decision variables are integer valued. In general, an MILP problem can be written as

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{y}}{\text{minimize}} && \mathbf{c}^T \mathbf{x} + \mathbf{d}^T \mathbf{y} \\ & \text{subject to} && \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} = \mathbf{b}, \\ & && \mathbf{x}, \mathbf{y} \geq 0, \\ & && \mathbf{x} \text{ integer,} \end{aligned}$$

where $\mathbf{A}, \mathbf{B}, \mathbf{b}$ are matrices or a vector with commensurate dimensions. If the integer constraint is relaxed, an MILP problem becomes an LP problem. The MILP problem has many applications (for example, [53]). For instance, 0–1 binary decision variables can be used to model logical “on–off” decisions. If x_1 and x_2 are both 0–1 binary decision variables, then the constraint $x_1 + x_2 = 1$ means that either $x_1 = 1$ or $x_2 = 1$ but not both. This modeling capability is not available with an LP, because LP decision variables can take fractional values. Another well-known example of MILP modeling is the traveling salesman problem, where a map of cities and pairwise distances between cities are given and the salesman has to make a shortest-distance tour

visiting each city exactly once. The traveling salesman problem has many important applications including circuit-board drilling and DNA sequencing. The MILP model of the traveling salesman problem cannot be relaxed to an LP model since the decision of whether or not a road is traversed is a binary one. The MILP problem is NP-hard, as it includes as a special case the 0–1 integer program. As a result, unless $P = NP$, it is impossible to find a polynomial-time algorithm to solve the MILP problem. This implies that the computational effort for solving MILP problems in general increases very rapidly as the size of the problem increases. For example, suppose that a basic computation requires 10^{-9} s to perform on a computer. On a graph with $|V| = 30$ nodes and $|E| = |V|(|V| - 1)/2$ edges, to solve the traveling salesman problem by enumeration requires $O(2^{|E|})$ basic computations, or about 2×10^{122} years. On the other hand, for the same graph, if instead the minimum cut problem is solved with a polynomial-time algorithm that requires $O(|V||E| + |V|^2 \log(|V|))$ basic computations, then the solution time is only about 17 ms. Nevertheless, solution algorithms for the MILP problem are well studied and well developed. They include, for instance, branch-and-bound methods and cutting-plane methods. Software implementations of MILP problem solution algorithms include, for example, CPLEX [69] and Gurobi [70].

written as a linear program. Hence, it can be solved efficiently to obtain an exact optimal solution. The optimal solution to (12) is feasible to problem (8) because $C_y(j, :) \Delta x = 1$ implies $C_y(j, :) \Delta x \neq 0$. Therefore, the optimal solution to (12) is an approximate solution to (8), with the former leading to an objective value that is greater than or equal to the true minimum of (8). Therefore, the ℓ_1 -relaxation approach provides an overestimate of the security index.

An alternative approach to handle the large-scale system computation difficulty is to develop specialized algorithms for particular instances of (8). For example, when the underlying application is power network state estimation and when the measurement system satisfies certain assumptions, such as the full measurement assumption to be described, problem (8) can be solved exactly in a time-efficient manner. The details of this result and an illustration with large-scale numerical examples when applied to electric power systems will be given in the following section.

Risk Analysis and Treatment for Electric Power Network

Power transmission networks are complex and spatially distributed systems. They are operated through SCADA systems, which represent the backbone IT and control infrastructure, as illustrated in Figure 7. SCADA systems collect data from remote terminal units (RTUs) installed in substations and relay aggregated measurements to the

central master station located at the control center. The technological limitations of legacy measurement equipment limits the sampling periods to the order of tens of seconds, thus the system is mainly observed at a quasi-static state.

SCADA systems for power networks are complemented by a set of application-specific software, usually called energy management systems (EMSs). EMSs enable state and measurement estimation and optimal operation under safety and reliability constraints by providing human operators with state-awareness and recommended control actions. In the past, the malfunction of EMS components, in particular the state estimator, has led to a large-scale blackout with severe economic consequences [58]. Furthermore, as discussed in [2], there are several vulnerabilities in the SCADA system architecture, including the direct tampering of RTUs, communication links between the RTUs and the control center, and the IT software and databases in the control center. Thus cybersecurity of SCADA and EMS in power networks is of major importance.

Given the relevance of power networks, in this part of the article the risk assessment method described in the previous section on static models will be specialized to the case where the plant is an electric power network. Focusing on the power network case enables the risk assessment to be performed in a computationally efficient manner. In addition, some of the risk treatment tools for power

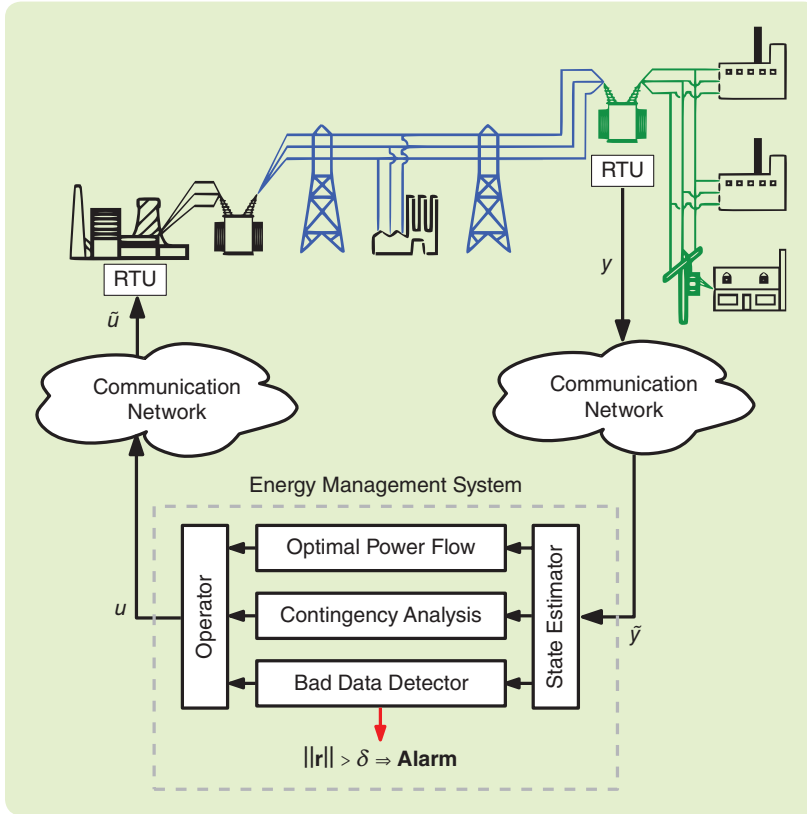


FIGURE 7 A schematic diagram of the electric power network and supervisory control and data acquisition (SCADA) system, adapted from [58]. Measurements taken from the remote terminal units (RTUs) are sent through the SCADA system to the control center. The received measurements are used by several energy management system applications that provide state awareness and control recommendations to human operators. The human operators decide the appropriate control actions and apply them to the power network through the SCADA system.

network applications will be highlighted. At the end of this section, a numerical case study with IEEE benchmark systems will illustrate the effectiveness of the risk assessment tools described in this section.

DC Power Flow Measurement Model

Assume that the electric power network has $n + 1$ buses and L transmission lines. The state of the network is determined by the complex voltages at the buses, whose magnitudes and phase angles are, respectively, denoted by V_i and x_i for $i = 0, 1, \dots, n$. In power networks, commonly considered measurements include line power flows, bus power injections, bus voltage magnitudes, and line current flow magnitudes. This section focuses only on active power flows on transmission lines and active power injections at buses, which are functions of the bus voltage magnitudes and phase angles. However, for the analysis of cyberphysical security, bad-data detection, and network observability, it is customary to describe the dependency of active power flows and injections through an approximate model called the dc-power flow model. By assuming that the voltage magnitudes V_i are all fixed to 1 p.u. (that

is, unity in the per unit system [42]), the dc power flow model depends only on the voltage phase angles. In this model, the transmission-line active power flow from bus i to bus j is

$$P_{ij} = \frac{x_{ij}}{X_{ij}}, \quad (13)$$

where $x_{ij} = x_i - x_j$ and $X_{ij} > 0$ is the reactance of the line between bus i and bus j . On the other hand, the active power injection at bus i is

$$P_i = \sum_{j \in N_i} P_{ij}, \quad (14)$$

where N_i is the set of all indices of the neighboring buses of bus i , excluding i .

Equations (13) and (14) give rise to a linear measurement model in matrix-vector form. Let x denote the n -vector of voltage phase angles on all buses except the reference bus. The reference bus is arbitrarily defined, with voltage phase angle fixed at zero. In addition, let y denote the vector of active power flow and active power injection measurements. Then, y and x are related by the equation

$$y = \begin{bmatrix} T_l D \mathcal{A}^T x \\ T_i \mathcal{A}_0 D \mathcal{A}^T x \end{bmatrix} =: C_y x. \quad (15)$$

In (15), the term $T_l D \mathcal{A}^T x$ corresponds to transmission line power flow measurements. On the other hand, the term

$T_i \mathcal{A}_0 D \mathcal{A}^T x$ corresponds to power injection measurements. The symbols in (15) are as follows: $\mathcal{A}_0 \in \mathbb{R}^{(n+1) \times L}$ is the incidence matrix of the network defined as

$$\mathcal{A}_0(i, l) = \begin{cases} 1 & \text{line } l \text{ starts from bus } i, \\ -1 & \text{line } l \text{ ends at bus } i, \\ 0 & \text{otherwise,} \end{cases} \quad \text{for each transmission line } l.$$

The directions of the lines in \mathcal{A}_0 are irrelevant to the application in this article. They can be fixed arbitrarily. Matrix \mathcal{A} is the truncated incidence matrix with the row of \mathcal{A}_0 corresponding to the reference bus removed. Matrix D is a diagonal matrix with the diagonal entries being the reciprocals of X_{ij} for all lines. Matrices T_l and T_i are stacked by the rows of identity matrices, and they indicate which line power flows and bus power injections are actually measured. The total number of rows of T_l and T_i is the total number of measurements, which is denoted by m . The matrix C_y is again referred to as the measurement matrix.

The measurement model in (15) has a network potential flow interpretation. A particular x corresponds to an

assignment of the bus voltage phase angles. The phase angle differences between two neighboring buses induce power flows along the connecting lines. The vector of induced line power flows is described by $D\mathcal{A}^T x$. At each bus, any difference between total incoming and total outgoing line power flows has to be balanced by an external power injection or extraction. The vector of external power injections at the buses is described by $\mathcal{A}_0 D\mathcal{A}^T x$. The matrices T_l and T_i allows the flexibility that all line power flows and bus power injections are not measured. The network potential flow interpretation is illustrated in Figure 8.

Security Index Problem Under dc Power Flow Model

In this section, unless otherwise noted, the security index problem in (8) should be interpreted with the measurement matrix C_y restricted to the form in (15). With the restriction of C_y , problem (8) has a network potential flow interpretation. The vector Δx can be considered as an assignment of fictitious voltage phase angles. The constraint in (8) states that either a particular line has a nonzero flow or a particular bus has a nonzero injection, depending on the meaning of measurement j . The objective is to minimize $\|C_y \Delta x\|_0 = \|T_l D\mathcal{A}^T \Delta x\|_0 + \|T_i \mathcal{A}_0 D\mathcal{A}^T \Delta x\|_0$, being the sum of the number of lines with nonzero measured flows and the number of buses with nonzero measured injections.

Even though C_y is specialized, (8) is still NP-hard [59]. However, (8) becomes solvable in polynomial time under certain additional assumptions. One such example is when T_l is an identity matrix and T_i is a zero matrix. In this case, minimizing the number of lines with nonzero phase angle differences is the only objective. It can be seen that an assignment of Δx using only two distinct values (say zero and one) is optimal. In fact, the above assumption can be generalized: the 0–1 assignment remains optimal if both T_l and T_i are identity matrices of appropriate dimensions [59]. In this case, all line power flows and bus injections are measured and the condition is referred to as the *full measurement assumption*. In summary, under the full measurement assumption, (8) is equivalent to

$$\begin{aligned} & \underset{\Delta x \in \{0,1\}^n}{\text{minimize}} && \|C_y \Delta x\|_0 \\ & \text{subject to} && C_y(j,:) \Delta x \neq 0. \end{aligned} \quad (16)$$

The only difference between (8) and (16) is that the decision vector of real numbers in (8) is replaced by the decision vector of 0–1 binary values in (16).

The 0–1 assignment of the entries of Δx in (16) leads to yet another graph interpretation. The binary choice of entries of Δx specifies a partitioning of the buses into two disjoint sets, a set with buses with fictitious voltage phase angles being zero and the complementary set. A line connecting two buses in two different sets is cut. The objective is to minimize the sum of the number of cut lines with line power flow meters and the number of buses that have injection meters and are incident to at least one cut line. The con-

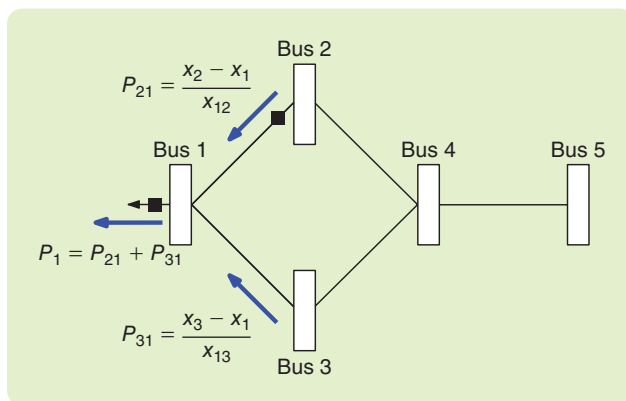


FIGURE 8 Potential flow interpretation of the line power flows and bus power injections. Buses 1–5 each have an assignment of voltage phase angle x_1, x_2, x_3, x_4 , and x_5 , respectively. The phase angle differences $x_2 - x_1$ and $x_3 - x_1$ each induce line power flows P_{21} on the line connecting bus 1 and bus 2 and P_{31} on the line connecting bus 1 and bus 3, respectively. The net injection of power flows into bus 1 is balanced by the external extraction P_1 , which is the bus injection at bus 1 with a negative value. All line power flows are stored in the vector $D\mathcal{A}^T x$. All power injections are stored in the vector $\mathcal{A}_0 D\mathcal{A}^T x$. The measurement vector z in (15) contains a subset of all line power flows and bus injections. The selection is achieved through T_l and T_i . In the figure, the black squares indicate the meters. Only the flows or injections with meters are measured.

straint can also be described as a particular line being cut [59]. As a result, (16) can be interpreted as a generalization of the standard minimum cut problem; see “Standard Minimum Cut Problem.” The only difference is that the standard minimum cut problem does not consider the cost associated with the number of buses incident to cut lines. The generalized minimum cut problem is illustrated in Figure 9.

The optimal solution to the generalized minimum cut problem can be interpreted in an alternative way: if all transmission lines in the cut are removed, then the remaining network contains at least two isolated subnetworks. This, in fact, implies that the measurement system with the compromised measurements removed is unobservable because it is impossible to deduce the phase angle information of one subnetwork from the information of another subnetwork [60]. In fact, in [26] it was shown that a set of compromised measurements is the optimal solution to the security index problem (8) if and only if the set is a sparsest critical tuple containing the target measurement j . A critical tuple is a set of measurements whose removal renders the remaining measurement system unobservable, but the removal of any strictly proper subset of the critical tuple would not lead to loss of observability. The interpretation of the security index problem optimal solution set as a critical tuple is illustrated in Figure 10. While this article focuses on interpreting the security index problem as a power network observability problem, the converse interpretation can be utilized to solve an observability analysis problem as a security index problem; see [61] for details.

Standard Minimum Cut Problem

A standard minimum cut problem (on an undirected graph) is an optimization problem whose instance is defined by an undirected graph with nonnegative edge weights and two distinct nodes called source and sink; see Figure S3 for an illustration.

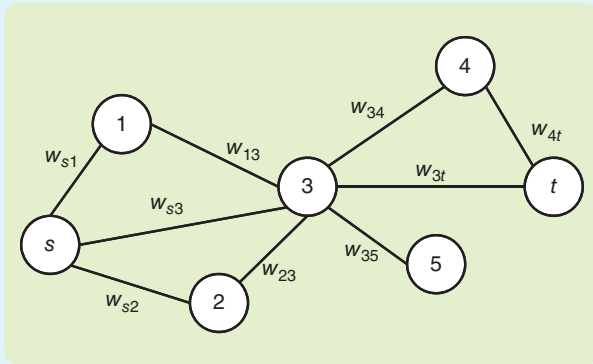


FIGURE S3 An instance of a minimum cut problem is defined on an undirected graph with nonnegative edge weights. The nodes s and t denote the source and sink, respectively. For an edge connecting node i and node j , the quantity $w_{ij} \geq 0$ is the corresponding edge weight.

The problem seeks a partition of the set of all nodes into two parts, with the source in one part and the sink in the other part, so that the sum of the weights of cut edges is minimized. An edge is cut if and only if it has one end node in the partition part including the source and the other end node in the partition part including the sink; see Figure S4 for an illustration. As an application, consider the network in Figure S3 as a power network s is the supply (generator) and t is the demand (load), respectively. If all

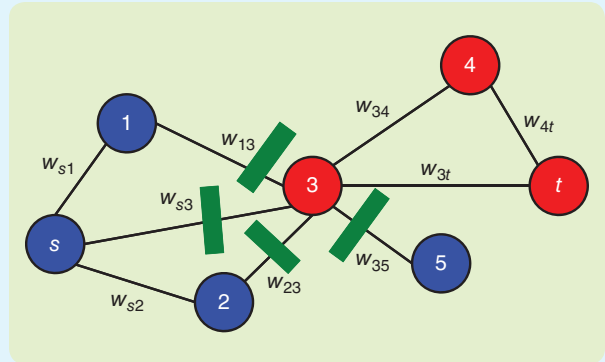


FIGURE S4 A feasible solution of the minimum cut problem is a partition of the nodes into two parts, color coded by blue and red in this illustration. A constraint of the minimum cut problem is that the source (node s) and the sink (node t) must be in two different parts. In this case, the source is blue and the sink is red. The partitioning of the nodes induces a cut of the edges. An edge is cut if and only if it connects nodes in different parts in the partition. In this illustration, only edges $\{1,3\}$, $\{s,3\}$, $\{2,3\}$, and $\{3,5\}$ are cut, as indicated by the green bars. The objective value corresponding to the particular node partition is the sum of the weights of cut edges, which is $w_{13} + w_{s3} + w_{23} + w_{35}$. The minimum cut problem seeks the partition of nodes separating s and t with the minimum total weights of cut edges.

edge weights are set to unity, then the solution to the minimum cut problem specifies the minimum number of transmission lines to break to cause the disruption of power supply.

Efficient solution algorithms are available to solve the standard minimum cut problem with computation effort proportional to a polynomial function of the size of the problem. The algorithms can be direct [62], or based on solving the dual maximum flow problem enabled by the max-flow, min-cut theorem [71].

The generalized minimum cut problem can be solved in polynomial time, due to its connection to the standard minimum cut problem, which is polynomial-time solvable [62]. The main result is that an instance of the generalized minimum cut problem is equivalent to an instance of the standard minimum cut problem on an auxiliary graph that is roughly three times as large as the original graph [59]. Additionally, the auxiliary graph can be constructed efficiently. As discussed in “Mixed Integer Linear Programming,” the minimum cut problem has much lower computational complexity than enumeration methods. Therefore, the generalized minimum cut formulation can potentially be applied to large problems, as illustrated in the case studies reported at the end of this section.

In summary, given an instance of the security index problem in (8) under the full measurement assumption or similar conditions found in [59], it can first be specialized by restricting the entries of Δx to either zero or one, as in

problem (16). The specialized problem can be viewed as an instance of the generalized minimum cut problem. Then, an instance of the standard minimum cut problem on an auxiliary graph can be set up and solved in polynomial time. The solution to the standard minimum cut problem can be used to construct the solution to the security index problem due to the equivalence mentioned above. Notice that, even without the full measurement assumption, the solution to problem (16) still serves as an overestimate of the true security index. In particular, the security index of a measurement must be small if the overestimate is small. This indicates the vulnerability of the measurement. The accuracy of this approach was demonstrated in [59].

Risk Treatment Approaches

One possible approach to decrease the risk of stealthy deception attacks is to encrypt the data and communication channels. Since a large part of today’s power grid equipment is

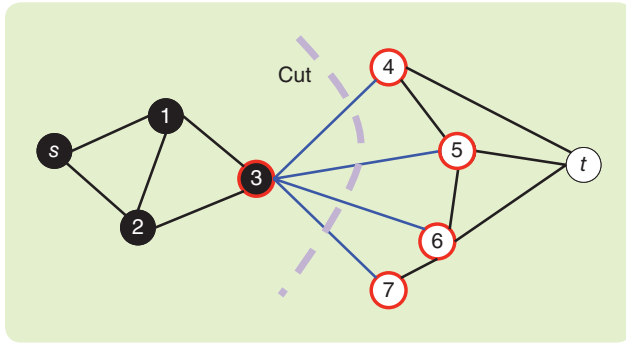


FIGURE 9 An illustration of the generalized minimum cut problem. An instance of the generalized minimum cut problem is defined by the graph with node weights p_i and edge weights w_{ij} , as well as the designation of the source node s and the sink node t . The node weights and the edge weights are nonnegative. They indicate the number of meters on the buses and lines, respectively. Any partition of the set of all nodes into two parts, with s being in one part and t being in the other part, corresponds to a feasible solution to the generalized minimum cut problem. The partition of the nodes induces a cut. An edge is cut if and only if it has an end node in one partition and the other end node in the other partition. The edges $\{3,4\}$, $\{3,5\}$, $\{3,6\}$, and $\{3,7\}$ are colored blue and are cut. The sum of the weights of the cut edges contribute to part of the objective value of the optimization. In addition, a node incident to a cut edge is also “in the cut.” In this example, nodes 3, 4, 5, 6, and 7 are in the cut and have red boundaries in the illustration. The sum of weights of the nodes in the cut contribute to the other part of the objective value. The total objective value associated with a node partition is the sum of all cut edge weights and all cut node weights.

old, data encryption can be costly to implement because of the corresponding update of the equipment. Therefore, the following question is of great importance to measurement data integrity. Given limited protection resources (the number of devices for data encryption), which measurements should be encrypted to maximize the benefits of the protection resources? The risk analysis outcome from computing the measurements’ security indices may be used to sort the measurements in terms of their vulnerability and identify those that should be protected. In fact, a variant of the security index problem can help provide an answer to the previous question, namely

$$\begin{aligned} & \underset{\Delta x \in \mathbb{R}^n}{\text{minimize}} && \|C_y \Delta x\|_0 \\ & \text{subject to} && C_y(j,:)\Delta x \neq 0, \\ & && C_y(P,:)\Delta x = 0, \end{aligned} \quad (17)$$

where P is the index set of the encrypted measurements, which cannot be attacked. By comparing the security indices for different index sets P , it is possible to evaluate the effect of different protection strategies and determine the best one to implement. For example, [28] considers a lexicographic optimization of some security metrics that are based on the security index computation related to (17).

In the case where it is impractical to encrypt all measurements, it becomes critical to detect and isolate the measurements that are under attack. Effective attack isolation

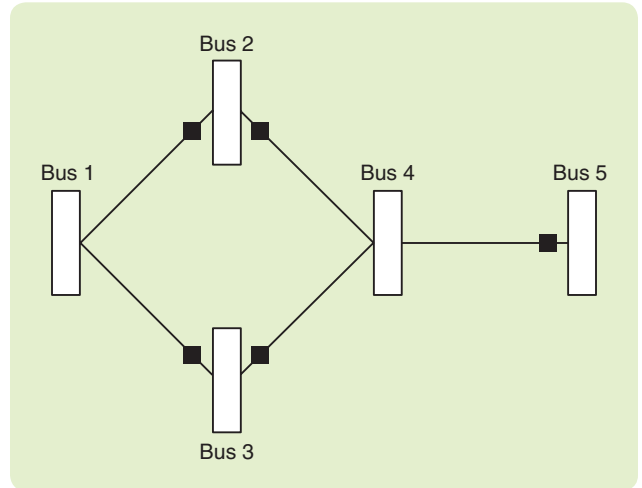


FIGURE 10 The optimal set of compromised measurements in a security index problem can be interpreted as a critical tuple of the measurement system. The security index problem with the line power flow measurement on line $\{4,5\}$ is solved with line $\{4,5\}$ being the optimal solution. Removing line $\{4,5\}$ renders bus 5 an isolated bus in the remaining network. It is then impossible to deduce the phase angle at bus 5 from the phase angle information of the remaining subnetwork and vice versa. The power flow measurement on line $\{4,5\}$ is a critical measurement (that is, critical one-tuple) and hence $\alpha = 1$. On the other hand, for the case where line $\{1,2\}$ is the target, the security index is two with three possible optimal sets of compromised lines: $\{\{1,2\},\{1,3\}\}$, $\{\{1,2\},\{3,4\}\}$, $\{\{1,2\},\{2,4\}\}$ but not $\{\{1,2\},\{4,5\}\}$ because the last set does not disconnect bus 1 from bus 2. All solution sets of the security index problem are critical pairs (that is, critical two-tuples, $\alpha = 2$), but $\{\{1,2\},\{4,5\}\}$ is not because removing the measurement on line $\{4,5\}$, which forms only a strictly proper subset of $\{\{1,2\},\{4,5\}\}$, still renders the measurement system unobservable.

enables the damage control (for example, removing attacked measurements for state estimation) to be performed in a timely fashion, that is, before the attack can lead to an incident with significant consequences. A distributed procedure for isolating the data attacks on power system transmission line power flow measurements is presented in [63], based on secure bus voltage magnitude measurements. The work in [18] develops a generalized likelihood ratio test to detect the presence of data attack, based on the assumption that the normal measurements follow a known Gaussian distribution. Mechanisms to detect data attacks based on known-secure phasor measurement unit PMU measurements and a known pattern of system states are presented in [29].

Security Index Problem Case Studies on the Benchmark Systems

As a numerical demonstration, the security index problems for all measurements in two benchmark systems are considered. The two benchmarks are the IEEE 14-bus and IEEE 118-bus systems [64]. In this case study, all line power flows and bus injections of the benchmark systems are measured. In other words, the full measurement assumption holds and

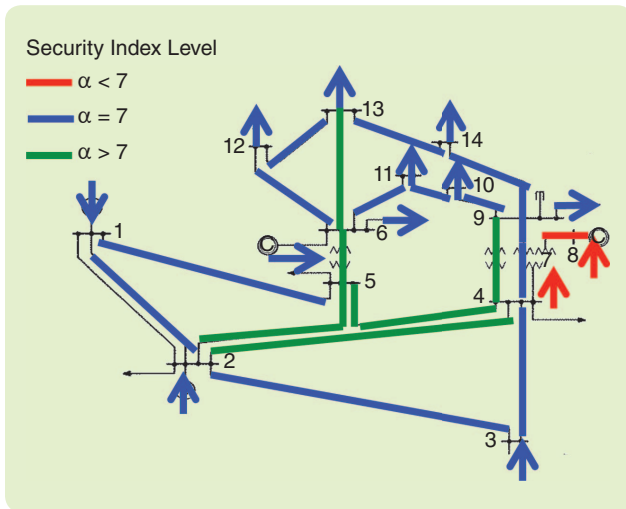


FIGURE 11 The IEEE 14-bus benchmark system with all measurements labeled different colors according to their resilience against stealthy data attack. The vulnerable measurements have small security indices (< 7) and are color coded red. The resilient measurements have large security indices (> 7) and are color coded green. The other measurements are color coded blue and their resilience lies somewhere in between. The efficient computation of security indices enables the rapid determination of the security weak points in the measurement system.

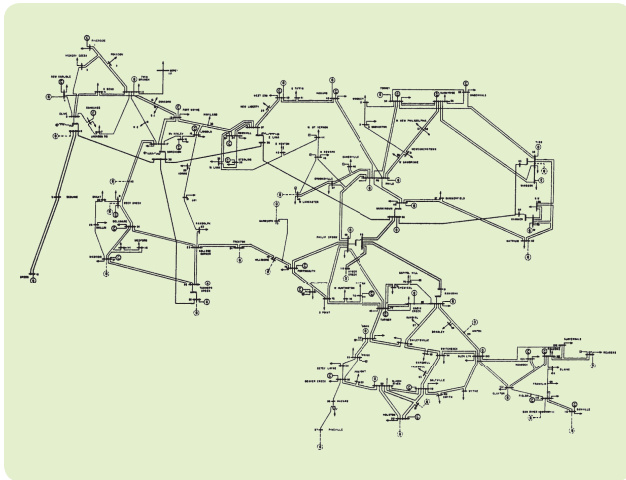


FIGURE 12 The IEEE 118-bus benchmark system [64].

the security index can be computed by solving (16). Figure 11 shows the IEEE 14-bus system with measurements color coded to indicate which ones are more vulnerable to stealthy data attack and which ones are resilient. The criterion to determine the vulnerability of the measurements is the security index. Measurements with small security indices (value lower than seven in Figure 11) are considered vulnerable and are color coded red. In particular, the injections at bus 7 and bus 8 as well as the line power flow between these two buses are vulnerable. This agrees with intuition, because bus 7 and bus 8 are on the boundary of the system with relatively little redundancy. On the other hand, the measurements located in the middle of the system (such as the

injection at bus 5) are considered resilient, with security indices much greater than seven. This also agrees with intuition: the measurements in the middle have great redundancy. The resilient measurements are color coded green in Figure 11. The rest of the measurements are color coded blue in Figure 11, with security index equal to seven. They are considered neither vulnerable nor resilient.

Next, the security index computation with the IEEE 118-bus system is considered. See Figure 12 for an illustration of this system. Figure 13 shows the sorted security indices for the measurements, computed using the minimum-cut-based procedure described in the previous subsection. The computation takes about 0.17 s on a personal laptop. On the other hand, when the big M method in (11) is used to compute the same security indices, the computation time is about 118 s. Figure 13 indicates that a significant number (about 40) of the measurements have relatively small security indices (values equal to four). In addition, there are a large number of measurements with security indices of seven. The efficient computation of the security indices is particularly relevant for real-world applications. Two motivations are presented below. First, the simplest realistic power network model contains at least thousands of buses (for example, the CAISO model contains 4000 buses). More realistic models in future applications are expected to grow in complexity. Thus, the time requirement for solving the security index problem can become excessive unless scalable computation procedures such as the minimum-cut-based one described in this article are available. Second, the security index computation presented so far is only for analysis purposes. In control design and synthesis situations, it is expected that the security index problem would need to be solved many times. For example, if there are m encryption devices to be employed, then it could take up to 2^m security index computations to determine the best configuration of encryption device deployment.

Risk Analysis for Dynamic Models

Having discussed methods for threat likelihood estimation in the static case, the dynamic case and methods for more general quantitative risk assessment are considered. The proposed risk assessment methods are not executed based on real-time data of the system. Instead, these methods are used offline, for a given configuration and corresponding model of the system, to assess the risk of different hypothetical attack scenarios. Consider the time interval $[0, N]$ and define the vectors $\mathbf{n} = [\eta_0^\top \dots \eta_N^\top]^\top$ and $\mathbf{a} = [a_0^\top \dots a_N^\top]^\top$, which capture the system states and attack signals over the time interval of interest. The state and residue trajectories are described by the following static mappings obtained (5)

$$\begin{aligned} \mathbf{n} &= \mathcal{O}\eta_0 + \mathcal{T}\mathbf{a}, \\ \mathbf{x} &= \mathbf{C}_x\mathbf{n}, \\ \mathbf{r} &= \mathbf{C}_r\mathbf{n} + \mathcal{D}_r\mathbf{a}, \end{aligned} \quad (18)$$

where \mathcal{O} describes the effect of the initial condition η_0 on the system's state trajectory, and \mathcal{T} is a lower triangular

block-Toeplitz matrix mapping the attack signals to the state trajectory; see [27] for details. For simplicity, η_0 is assumed to be zero. These mappings can be used for risk assessment of deception attacks on networked control systems, as described in the remainder of this section.

Maximum-Impact, Minimum-Resource Attacks

For illustration purposes, the quantitative risk assessment methods in the previous section considered only the adversary resources. In this subsection, full risk assessment is performed on dynamical systems by simultaneously considering impact and resources. Recall that the adversary aims to perturb the networked control system operation and drive the system to an unsafe set. To better clarify the proposed approach, suppose that the safe set is defined as

$$\mathcal{S}_x \triangleq \{\mathbf{x} : \|\mathbf{x}\|_p < \delta\},$$

where $\|\mathbf{x}\|_p$ denotes the ℓ_p -norm of \mathbf{x} for $p \geq 1$, and \mathbf{x} depends on the attack signal \mathbf{a} as described in (18). Therefore, the attack impact during the time interval $[0, N]$ can be characterized as the perturbation of the state trajectory \mathbf{x} due to the attack quantified by $\|\mathbf{x}\|_p$. Notice that the proposed framework can be straightforwardly applied to safe sets that consider linear transformations of the state trajectory.

To illustrate the resources required for a given attack signal, suppose that Δu_k and Δy_k are scalars and recall the attack vector at time k , $\mathbf{a}_k = [\Delta u_k \ \Delta y_k]^\top$. Consider an attack having Δu_k equal to zero for all k , while Δy_k is equal to zero at all times except for $k = 0$. Since Δy_0 is nonzero, the adversary must have access to one communication channel to inject the false data, namely the measurement channel. The attack with $\Delta y_k \neq 0$ only at $k = 0$ requires as many resources as the one having $\Delta y_k \neq 0$ for $k > 2$, thus corrupting the measurement signal at other times does not require additional adversary resources. On the other hand, corrupting also the control signal such that $\Delta u_k \neq 0$ at a given time k requires an additional resource, namely the access to the control signal. More generally, the attack vector can be rewritten as $\mathbf{a}_k = [a_{(1),k}, \dots, a_{(n_a),k}]^\top$, where $a_{(i),k} \in \mathbb{R}$ denotes the corrupted data introduced in the i th adversary resource at time k . Denoting the attack signal at the i th resource by $\mathbf{a}_{(i)} = [a_{(i),0}, \dots, a_{(i),N}]^\top$, the i th resource is used during the attack if $\|\mathbf{a}_{(i)}\|_p$ is nonzero. Defining the vector

$$h_p(\mathbf{a}) \triangleq [\|\mathbf{a}_{(1)}\|_p \ \dots \ \|\mathbf{a}_{(n_a)}\|_p]^\top,$$

the number of resources used in a given attack corresponds to the number of nonzero elements of $h_p(\mathbf{a})$, which is denoted as $\|h_p(\mathbf{a})\|_0$.

The attack impact and resources are jointly considered in the multiobjective optimization problem [65]

$$\begin{aligned} & \underset{\mathbf{a}}{\text{maximize}} \quad [\|\mathbf{x}\|_p, -\|h_p(\mathbf{a})\|_0]^\top \\ & \text{subject to} \quad \|\mathbf{r}\|_q \leq \delta, \end{aligned} \quad (19)$$

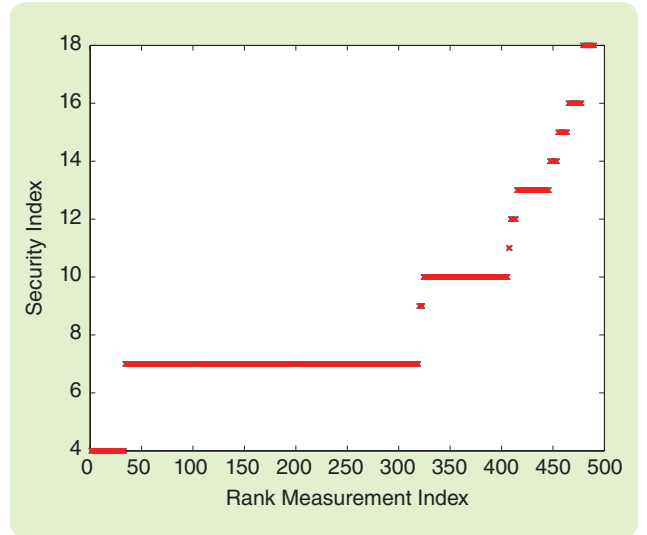


FIGURE 13 Sorted security indices for the fully measured 118-bus system. The minimum-cut-based procedure computes the indices in about 0.17 s, while the big M method in (11) requires 118 s to obtain the same result. This figure indicates that a significant number (about 40) of the measurements have relatively small security indices (values equal to four). In addition, there are a large number of measurements with security indices of seven. The efficient computation of the security indices allows the network operator to quickly assess the security of the power network and to identify the weak points.

where \mathbf{x} and \mathbf{r} are determined by \mathbf{a} according to (18). The multiobjective optimization problem (19) can be interpreted as computing the attack signal \mathbf{a} that simultaneously maximizes the impact $\|\mathbf{x}\|_p$ and minimizes the resources $\|h_p(\mathbf{a})\|_0$, while remaining stealthy by ensuring $\|\mathbf{r}\|_q \leq \delta$.

The tradeoff analysis in the multiobjective optimization problem can be performed by studying the Pareto solutions [65]. These solutions can be obtained through several techniques, for instance, the bounded objective function method in which all but one of the objectives are posed as constraints, thus obtaining a scalar-valued objective function. Applying this method to (19) and constraining $\|h_p(\mathbf{a})\|_0$ yields

$$\begin{aligned} & \underset{\mathbf{a}}{\text{maximize}} \quad \|\mathbf{x}\|_p \\ & \text{subject to} \quad \|\mathbf{r}\|_q \leq \delta, \\ & \quad \quad \quad \|h_p(\mathbf{a})\|_0 < \epsilon, \\ & \quad \quad \quad \mathbf{n} = \mathcal{O}\eta_0 + \mathcal{T}\mathbf{a}, \\ & \quad \quad \quad \mathbf{x} = \mathcal{C}_x\mathbf{n}, \\ & \quad \quad \quad \mathbf{r} = \mathcal{C}_r\mathbf{n} + \mathcal{D}_r\mathbf{a}, \end{aligned} \quad (20)$$

which can be interpreted as a maximum-impact resource-constrained attack policy.

The Pareto frontier that characterizes the optimal tradeoff manifold can be obtained by iteratively solving (20) for $\epsilon \in \{1, \dots, n_a\}$. For a fixed ϵ and $p = q = \infty$, the optimization problem (20) can be formulated as an MILP, while the problem with parameters $\epsilon = +\infty$ and $p = q = 2$ reduces to a generalized eigenvalue problem; see [27] for a detailed discussion.

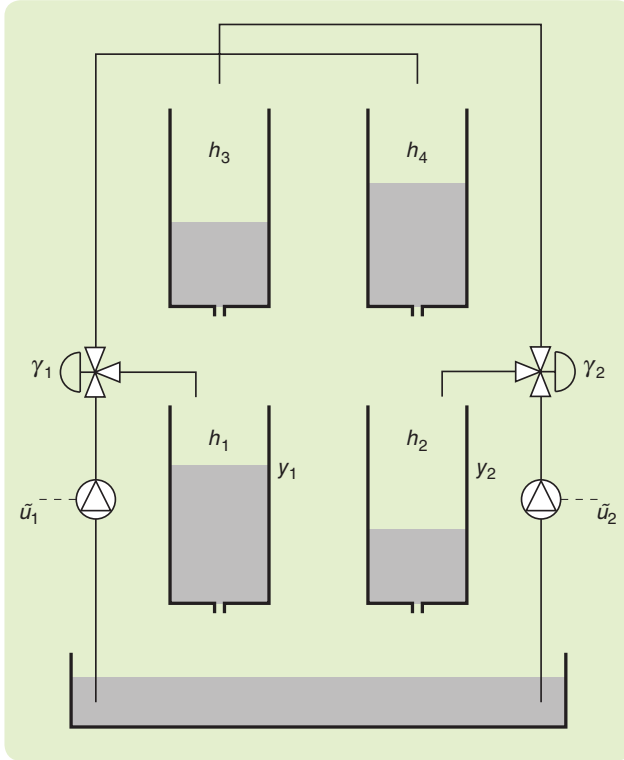


FIGURE 14 A schematic of the quadruple-tank process [66]. The water pump to the left is controlled by the signal \tilde{u}_1 , while the pump to the right is controlled through \tilde{u}_2 . The left pump forces water into tanks 1 and 4, where the fraction of water flowing into each tank is determined by the valve position γ_1 . Similarly, the fraction of water pumped to tanks 2 and 3 by the right-hand-side pump depends on γ_2 .

Case Study: Quadruple-Tank Process

Risk analysis and risk treatment approaches for dynamic control systems are illustrated for a particular networked control system in this section. The physical plant consists of the quadruple-tank process (QTP) [66], depicted in Figure 14, while the networked control system architecture is depicted in Figure 15.

The plant model is

$$\begin{aligned} \frac{dh_1}{dt} &= -\frac{a_1}{A_1}\sqrt{2gh_1} + \frac{a_3}{A_1}\sqrt{2gh_3} + \frac{\gamma_1 k_1}{A_1}u_1, \\ \frac{dh_2}{dt} &= -\frac{a_2}{A_2}\sqrt{2gh_2} + \frac{a_4}{A_2}\sqrt{2gh_4} + \frac{\gamma_2 k_2}{A_2}u_2, \\ \frac{dh_3}{dt} &= -\frac{a_3}{A_3}\sqrt{2gh_3} + \frac{(1-\gamma_2)k_2}{A_3}u_2, \\ \frac{dh_4}{dt} &= -\frac{a_4}{A_4}\sqrt{2gh_4} + \frac{(1-\gamma_1)k_1}{A_4}u_1, \end{aligned}$$

where $h_i \in [0, 30]$ are the heights of water in each tank, A_i is the cross-sectional area of each tank, a_i is the cross-sectional area of each tank's outlet, k_i the pump constants, γ_i the flow ratios, and g is gravitational acceleration. The nonlinear plant model is linearized for a given operating point and the state of the linearized plant model x corresponds to the water level deviations from the operating

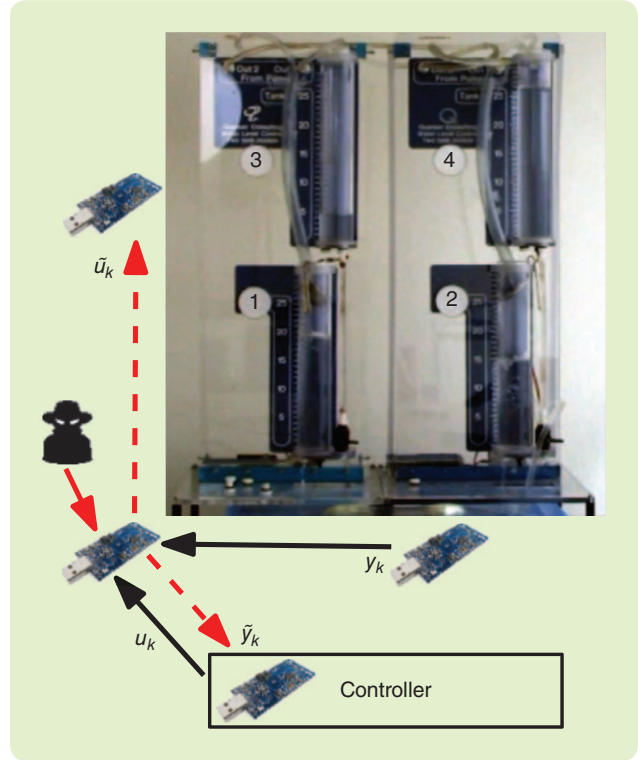


FIGURE 15 A schematic diagram of the test bed with the quadruple-tank process and a multihop communication network. The controller, sensors, and actuators communicate using a wireless network that has one relay node. The adversary is able to compromise the relay node and therefore has access to the control and measurement channels.

point. Moreover, given the range of the water levels and the operating point, the safe set is considered to be $\mathcal{S}_x = \{x \in \mathbb{R}^{n_x}: \|x\|_\infty \leq 5\}$.

The process is controlled using a centralized linear-quadratic-Gaussian (LQG) controller with integral action running on a remote computer and a wireless network is used for the communications. A Kalman filter based anomaly detector is also running on the remote computer and alarms are triggered according to (4). The measurements of the lower tanks' water levels, y_1 and y_2 , are sent to the controller and anomaly detector through the communication network. Given the received measurements \tilde{y}_1 and \tilde{y}_2 , the LQG controller computes the pump control signals, u_1 and u_2 , which are sent to the water pumps through the wireless network. In the present attack scenario, the adversary can corrupt the data transmitted through the wireless network, namely u_1 , u_2 , y_1 , and y_2 .

Risk Analysis for Stealthy Deception Attacks

For the time interval $[0, 50]$, the maximum-impact, minimum-resource attacks were computed for the process in minimum and nonminimum phase settings by choosing $p = q = 2$ and iteratively solving (20) with respect to ϵ . The respective impacts correspond to the energy of the state signal x for value of ϵ , and are presented in Table 1, while

the risk is depicted by the risk matrix plot in Figure 16(a) (recall Figure 5).

As expected, the nonminimum phase system is less resilient than the minimum-phase one. In both settings, the attack impact can be made arbitrarily large by corrupting three or more channels, and thus the adversary can drive the state out of the safe set while remaining stealthy. The results indicate that the threats compromising three or more channels have high risk and should therefore be analyzed in more detail. The risk of such threats can be mitigated by protecting the data channels, which is shown in the next subsection.

For illustrative purposes, the maximum-impact attack signal for the nonminimum phase system with $\epsilon = 2$, $\delta = 0.15$, and $p = q = 2$ is presented in Figure 17(a). For the parameters $\epsilon = 2$, $\delta = 0.025$, and $p = q = \infty$, the maximum-impact attack signal shown in Figure 17(b) was computed by solving an MILP; see [27]. In both cases the optimal attack corrupts both actuator channels and ensures $\|\mathbf{r}\|_p \leq \delta$, while maximizing the attack impact. Although the impact results in Table 1 do not quantify the impact according to the safe set $\mathcal{S}_x = \{x \in \mathbb{R}^{n_x} : \|x\|_\infty \leq 5\}$, the state trajectory does leave the safe set in both cases.

TABLE 1 Risk analysis results for the quadruple-tank process. Each entry corresponds to the maximum impact $\|\mathbf{x}\|_p$ for a given number of corrupted channels, computed through (20), with $p = q = 2$ and $\delta = 0.15$.

	Number of Compromised Channels			
	4	3	2	1
Minimum phase	∞	∞	140.39	1.15
Nonminimum phase	∞	∞	689.43	2.8

The attack signals illustrated in Figure 17 are related to the zero dynamics of the QTP system. The zero-dynamics attack signal and other scenarios were analyzed and performed in an experimental test bed of the QTP by [13]. Videos of the experiments are available [67].

Risk Treatment Approaches

The risk analysis identifies the data channels that, when corrupted, may lead to a large impact on the system. The subsequent step in the risk management framework is the risk

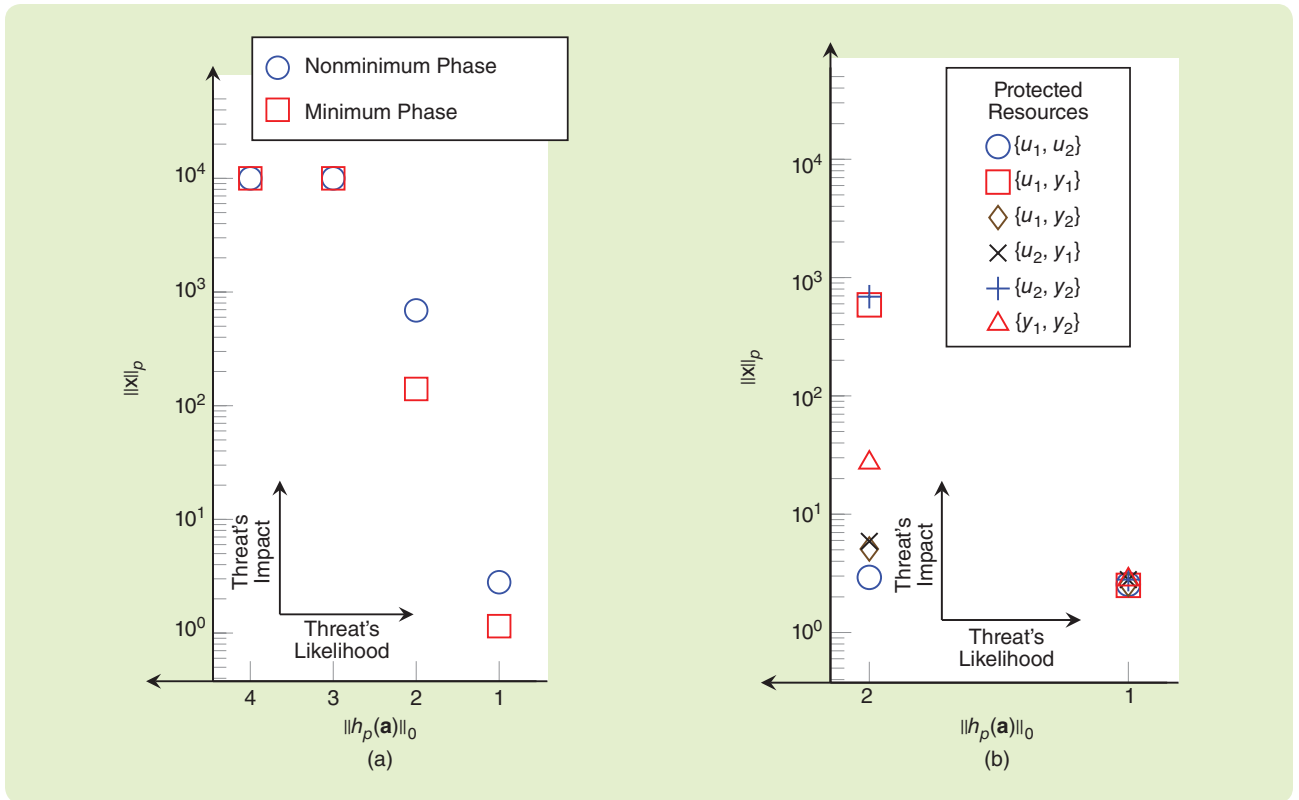


FIGURE 16 The risk matrix plot for the quadruple-tank process (QTP) (a) without protection and (b) for the nonminimum phase case when different pairs of resources are protected. The threat's likelihood is taken as a decreasing function of the number of compromised data channels, $\|h_p(\mathbf{a})\|_0$, and corresponds to the x-axis. The threat's impact on the y axis is the ℓ_p -norm of the state trajectory, $\|\mathbf{x}\|_p$. In (a), the risk analysis results for the minimum phase system (square) and nonminimum phase (circle) from Table 1 are depicted and qualitatively classified. The figure in (b) indicates that, when pairs of resources can be protected in the nonminimum phase process, the most effective choice for risk treatment is to protect both actuator channels, $\{u_1, u_2\}$.

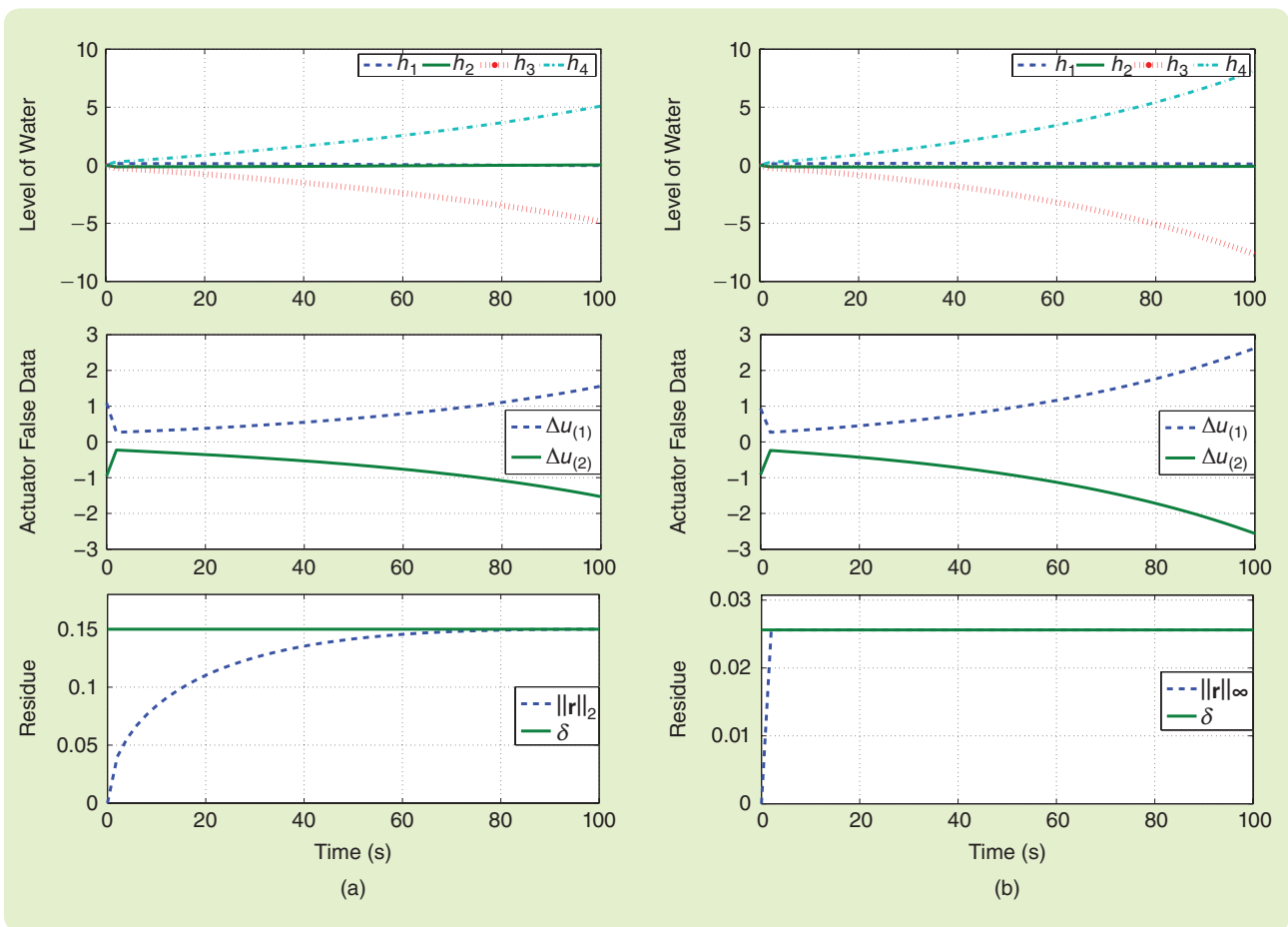


FIGURE 17 Simulation results for the maximum-impact attack signal. The attack signal is computed by solving the multiobjective problem (20) with $\epsilon = 2$ for the nonminimum phase system using the parameters (a) $p = q = 2$, $\delta = 0.15$ and (b) $p = q = \infty$, $\delta = 0.025$. In (a) and (b), the top plot depicts the water level change in each tank, the middle plot illustrates the false-data signal injected in the first and second actuators, $\Delta u_{(1)}$ and $\Delta u_{(2)}$, respectively, while the bottom plot shows the running energy and peak of the residue signal.

treatment stage, in which actions reducing the risk are chosen and implemented. A common approach to decrease the risk of threats is to deploy protective resources such as encryption, thus preventing the attacks from occurring. To assess the effectiveness of protecting a given set of data channels P the optimization problem (20) may be modified to

$$\begin{aligned}
 & \underset{\mathbf{a}}{\text{maximize}} && \|\mathbf{x}\|_p \\
 & \text{subject to} && \|\mathbf{r}\|_p \leq \delta, \\
 & && \|h_p(\mathbf{a})\|_0 < \epsilon, \\
 & && (18), \\
 & && \mathbf{a}_{(i)} = 0, \text{ for all } i \in P.
 \end{aligned} \tag{21}$$

The QTP example is now considered to illustrate the risk treatment step using channel encryption. The preventive action under study is the encryption of one pair of data channels so that the risk is minimized. The optimization problem (21) is solved for each pair of data channels, and the corresponding risk matrices plots are depicted in Figure 16(b).

From the results in Figure 16(b), it is evident that the pair of actuators $\{u_1, u_2\}$ should be protected to minimize the risk. Moreover, recalling the original risk matrix plot in

Figure 16(a), observe that the impact when two channels are corrupted is substantially decreased by protecting $\{u_1, u_2\}$. This protection choice is expected since the adversary can no longer inject an attack exciting the unstable zero dynamics of the system when both actuators are protected. Furthermore, since the resources accessible to the adversary are y_1 and y_2 , the adversary cannot have a direct impact on the physical system but instead needs to affect the system through the feedback controller by corrupting the measurement signals.

Methods other than encryption have been proposed in the literature to reduce the risk of threats. Concerning replay attacks, [36] proposes the use of a hypothesis test as the anomaly detector and the injection of random, zero-mean Gaussian noise with an optimally designed covariance in the control input channels. The injected noise increases the performance of the hypothesis test since the noise statistics are assumed to be unknown to the adversary. Similarly, [30] proposes the insertion of uncertainty in the adversary's model knowledge by modifying the system dynamics and control and output channels. The effects of such actions on zero-dynamics attacks are also characterized in detail.

Unlike other IT systems where cybersecurity mainly involves the protection of data, cyberattacks on networked control systems may influence physical processes through feedback actuation.

CONCLUSIONS AND FUTURE WORK

The pervasive use of IT infrastructures supporting the operation of networked control systems has introduced vulnerabilities into these systems, raising numerous challenges regarding the cyberphysical security of networked control systems. The specific nature of these threats and the coupling between the cyber and physical realms of the system requires the development of new paradigms and frameworks to study and tackle security-related problems.

This article described a cybersecurity problem in networked control systems, covering some of the key aspects such as the networked control-system architecture, the adversary model, and the defense methodology. The networked control-system architecture consists of the physical plant, the feedback controller colocated with the anomaly detector, and the communication network, through which the measurement and control data are sent. The adversary is modeled by a resource-constrained policy with limited model knowledge, disruption resources, and disclosure resources. Moreover, the attack policy is shaped according to the adversary's intent: to drive the state of the physical plant to an unsafe region while remaining undetected by the anomaly detector. The defense methodology is based on the risk management framework, where the concept of risk is defined as a function of the threat's likelihood and the threat's impact to the system. The risk management cycle continuously minimizes the risk of threats by performing risk analysis, risk treatment, and risk monitoring.

Recent quantitative methods developed for risk analysis are also presented, which are of major importance to enhance the cyberphysical security of networked control systems. The problem of quantifying the likelihood of threats for static systems is discussed, and computationally efficient methods are described and illustrated for large-scale electric power networks. Possible risk treatment approaches proposed in the literature are also mentioned.

Finally, the full risk analysis problem for dynamic systems simultaneously considering the attack impact and likelihood is described. The quantitative risk analysis methods are formulated and illustrated on the wireless QTP test bed. Additionally, the effectiveness of protection-based risk treatment schemes is evaluated in terms of their effect on the risk of threats, and alternative approaches proposed in the literature are summarized.

Possible extensions to the work described in this article include the study of efficient tools for risk assessment of

dynamic systems and the analysis of attack and defense policies under noise and uncertain communication channels, among several other interesting directions.

ACKNOWLEDGMENTS

The research leading to these results received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement 608224, the Swedish Research Council under grants 2009-4565 and 2013-5523, and the Knut and Alice Wallenberg Foundation.

AUTHOR INFORMATION

André Teixeira (andretei@kth.se) is a Ph.D. candidate in automatic control at KTH Royal Institute of Technology, Stockholm, Sweden. He received the M.Sc. degree in electrical and computer engineering in 2009 from the Faculdade de Engenharia da Universidade do Porto, Portugal. He was a finalist for the NecSys 2012 Best Student Paper Award and one of his publications is listed in ACM Computing Review's Notable Computing Books and Articles of 2012. His main research interests include secure control, distributed fault detection and isolation, distributed optimization, power systems, and multiagent systems. He can be contacted at the Department of Automatic Control, KTH School of Electrical Engineering, Osqualdas väg 10, plan 6, SE-100 44 Stockholm, Sweden.

Kin Cheong Sou received a Ph.D. degree in electrical engineering and computer science at the Massachusetts Institute of Technology in 2008. From 2008 to 2010 he was a postdoctoral researcher at Lund University, Sweden. From 2010 to 2013 he was a postdoctoral researcher at KTH Royal Institute of Technology, Stockholm, Sweden. Since 2013 he has been an assistant professor with the Department of Mathematical Sciences at Chalmers University of Technology and the University of Gothenburg, Sweden. His research interests include power-system cybersecurity analysis, environment-aware buildings and communities, convex and nonconvex optimization, and model reduction for dynamical systems.

Henrik Sandberg received the M.Sc. degree in engineering physics and the Ph.D. degree in automatic control from Lund University, Sweden, in 1999 and 2004, respectively. He is an associate professor with the Automatic Control Laboratory, KTH Royal Institute of Technology, Stockholm, Sweden. From 2005 to 2007, he was a post-doctoral scholar with the California Institute of Technology, Pasadena. In 2013, he was a visiting scholar at the Laboratory for

Information and Decision Systems at the Massachusetts Institute of Technology, Cambridge. He has also held visiting appointments with the Australian National University and the University of Melbourne, Australia. His current research interests include secure networked control, power systems, model reduction, and fundamental limitations in control. He was a recipient of the Best Student Paper Award from the IEEE Conference on Decision and Control in 2004 and an Ingvar Carlsson Award from the Swedish Foundation for Strategic Research in 2007. He is currently an associate editor of *Automatica*.

Karl Henrik Johansson is director of the KTH ACCESS Linnaeus Centre and professor at the School of Electrical Engineering, Royal Institute of Technology, Sweden. He is a Wallenberg Scholar and has held a six-year senior researcher position with the Swedish Research Council. He is director of the Stockholm Strategic Research Area ICT The Next Generation. He received the M.Sc. and Ph.D. degrees in electrical engineering from Lund University. He has held visiting positions at the University of California, Berkeley (1998–2000) and the California Institute of Technology (2006–2007). His research interests are in networked control systems; hybrid and embedded systems; and applications in transportation, energy, and automation systems. He has been a member of the IEEE Control Systems Society Board of Governors and chair of the IFAC Technical Committee on Networked Systems. He has been on the editorial boards of several journals, including *Automatica*, *IEEE Transactions on Automatic Control*, and *IET Control Theory and Applications*. He is currently on the editorial board of *IEEE Transactions on Control of Networked Systems* and the *European Journal of Control*. He has been guest editor for special issues, including the 2011 issue on wireless sensor and actuator networks in *IEEE Transactions on Automatic Control*. He was the general chair of the ACM/IEEE Cyberphysical Systems Week 2010 in Stockholm and IPC chair of many conferences. He has served on the executive committees of several European research projects in the area of networked embedded systems. In 2009, he received the Best Paper Award of the IEEE International Conference on Mobile Ad-hoc and Sensor Systems. In 2009, he was also appointed as a Wallenberg Scholar as one of the first ten scholars from all sciences, by the Knut and Alice Wallenberg Foundation. He was awarded an Individual Grant for the Advancement of Research Leaders from the Swedish Foundation for Strategic Research in 2005. He received the triennial Young Author Prize from IFAC in 1996 and the Peccei Award from the International Institute of System Analysis, Austria, in 1993. He received Young Researcher Awards from Scania in 1996 and from Ericsson in 1998 and 1999. He is a Fellow of the IEEE.

REFERENCES

[1] J. Hespanha, P. Naghshtabrizi, and Y. Xu, "A survey of recent results in networked control systems," *Proc. IEEE*, vol. 95, no. 1, pp. 138–162, Jan. 2007.
 [2] 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. Resilient Control Systems*, Idaho Falls, ID, Aug. 2009, pp. 31–35.

[3] S. Gorman, "Electricity grid in U.S. penetrated by spies," *Wall Street J.*, p. A1, Apr. 8, 2009.
 [4] N. Falliere, L. Murchu, and E. Chien. (2011, Feb.). W32.Stuxnet dossier. *Symantec*. [Online]. Available: www.symantec.com/content/en/us/enterprise/media/security_response/whitepapers/w32_stuxnet_dossier.pdf
 [5] T. Rid, "Cyber war will not take place," *J. Strategic Studies*, vol. 35, no. 1, pp. 5–32, 2011.
 [6] K. Zhou, J. C. Doyle, and K. Glover, *Robust and Optimal Control*. Upper Saddle River, NJ: Prentice-Hall, Inc., 1996.
 [7] I. Hwang, S. Kim, Y. Kim, and C. E. Seah, "A survey of fault detection, isolation, and reconfiguration methods," *IEEE Trans. Control Syst. Technol.*, vol. 18, no. 3, pp. 636–653, May 2010.
 [8] A. Teixeira, G. Dán, H. Sandberg, and K. H. Johansson, "Cyber security study of a SCADA energy management system: Stealthy deception attacks on the state estimator," in *Proc. 18th IFAC World Congr.*, Milano, Italy, Aug.–Sept. 2011, pp. 11271–11277.
 [9] R. Smith, "A decoupled feedback structure for covertly appropriating networked control systems," in *Proc. 18th IFAC World Congr.*, Milano, Italy, Aug.–Sept. 2011, pp. 90–95.
 [10] F. Pasqualetti, F. Dorfler, and F. Bullo, "Attack detection and identification in cyber-physical systems," *IEEE Trans. Autom. Contr.*, vol. 58, no. 11, pp. 2715–2729, Nov. 2013.
 [11] S. Sridhar, A. Hahn, and M. Govindarasu, "Cyber-physical system security for the electric power grid," *Proc. IEEE*, vol. 100, no. 1, pp. 210–224, 2012.
 [12] A. Cárdenas, S. Amin, and S. Sastry, "Research challenges for the security of control systems," in *Proc. 3rd USENIX Workshop Hot Topics Security*, San Jose, CA, July 2008, pp. 6:1–6:6.
 [13] A. Teixeira, I. Shames, H. Sandberg, and K. H. Johansson, "A secure control framework for resource-limited adversaries," *Automatica*, 2014. Available: dx.doi.org/10.1016/j.automatica.2014.10.067
 [14] A. Cárdenas, S. Amin, Z. Lin, Y. Huang, C. Huang, and S. Sastry, "Attacks against process control systems: Risk assessment, detection, and response," in *Proc. 6th ACM Symp. Information, Computer Communications Security*, New York, 2011, pp. 355–366.
 [15] H. Sandberg, A. Teixeira, and K. H. Johansson, "On security indices for state estimators in power networks," in *Proc. Preprints 1st Workshop Secure Control Systems, CPSWEEK*, Stockholm, Sweden, Apr. 2010.
 [16] L. Xie, Y. Mo, and B. Sinopoli, "False data injection attacks in electricity markets," in *Proc. 1st IEEE Int. Conf. Smart Grid Communications*, Gaithersburg, MD, Oct. 2010, pp. 226–231.
 [17] A. Teixeira, H. Sandberg, G. Dán, and K. H. Johansson, "Optimal power flow: Closing the loop over corrupted data," in *Proc. American Control Conf.*, Montreal, QC, Canada, June 2012, pp. 3534–3540.
 [18] O. Kosut, L. Jia, R. J. Thomas, and L. Tong, "Malicious data attacks on the smart grid," *IEEE Trans. Smart Grid*, vol. 2, no. 4, pp. 645–658, 2011.
 [19] Y. Mo and B. Sinopoli, "Secure control against replay attack," in *Proc. 47th Annu. Allerton Conf. Communication, Control, Computing*, Allerton, IL, Oct. 2009, pp. 911–918.
 [20] A. Gupta, C. Langbort, and T. Başar, "Optimal control in the presence of an intelligent jammer with limited actions," in *Proc. 49th IEEE Conf. Decision Control*, Atlanta, GA, Dec. 2010, pp. 1096–1101.
 [21] H. Fawzi, P. Tabuada, and S. Diggavi, "Security for control systems under sensor and actuator attacks," in *Proc. 51st IEEE Conf. Decision Control*, Maui, Hawaii, Dec. 2012, pp. 3412–3417.
 [22] S. Sundaram, S. Revzen, and G. Pappas, "A control-theoretic approach to disseminating values and overcoming malicious links in wireless networks," *Automatica*, vol. 48, no. 11, pp. 2894–2901, 2012.
 [23] F. Pasqualetti, A. Bicchi, and F. Bullo, "Consensus computation in unreliable networks: A system theoretic approach," *IEEE Trans. Autom. Contr.*, vol. 57, no. 1, pp. 90–104, 2012.
 [24] A. Khamis, B. Touri, and T. Başar, "Consensus in the presence of an adversary," in *Proc. 3rd IFAC Workshop Estimation Control Networked Systems*, Santa Barbara, CA, Sept. 2012, pp. 276–281.
 [25] Q. Zhu and T. Başar, "A dynamic game-theoretic approach to resilient control system design for cascading failures," in *Proc. 1st Int. Conf. High Confidence Networked Systems*, New York, 2012, pp. 41–46.
 [26] K. C. Sou, H. Sandberg, and K. Johansson, "On the exact solution to a smart grid cyber-security analysis problem," *IEEE Trans. Smart Grid*, vol. 4, no. 2, pp. 856–865, 2013.

- [27] A. Teixeira, K. C. Sou, H. Sandberg, and K. H. Johansson, "Quantifying cyber-security for networked control systems," in *Control of Cyber-Physical Systems* (Lecture Notes in Control and Information Sciences, no. 449), D. C. Tarraf, Ed. Switzerland: Springer Int. Publishing, 2013, pp. 123–142.
- [28] O. Vukovic, K. C. Sou, G. Dan, and H. Sandberg, "Network-aware mitigation of data integrity attacks on power system state estimation," *IEEE J. Select. Areas Commun.*, vol. 30, no. 6, pp. 1108–1118, 2012.
- [29] A. Giani, E. Bitar, M. Garcia, M. McQueen, P. Khargonekar, and K. Poola, "Smart grid data integrity attacks," *IEEE Trans. Smart Grid*, vol. 4, no. 3, pp. 1244–1253, Sept. 2013.
- [30] A. Teixeira, I. Shames, H. Sandberg, and K. H. Johansson, "Revealing stealthy attacks in control systems," in *Proc. 50th Annu. Allerton Conf. Communication, Control, Computing*, Allerton, IL, Oct. 2012, pp. 1806–1813.
- [31] S. X. Ding, *Model-Based Fault Diagnosis Techniques: Design Schemes*. New York: Springer Verlag, 2008.
- [32] H. Fawzi, P. Tabuada, and S. Diggavi, "Secure estimation and control for cyber-physical systems under adversarial attacks," *IEEE Trans. Autom. Contr.*, vol. 50, no. 6, pp. 1454–1467, 2014.
- [33] Y. Shoukry and P. Tabuada. (2013, Sept.). Event-triggered state observers for sparse sensor noise/attacks. *ArXiv e-Prints*. [Online]. Available: arxiv.org/abs/1309.3511
- [34] P. M. Frank and X. Ding, "Survey of robust residual generation and evaluation methods in observer-based fault detection systems," *J. Process Control*, vol. 7, no. 6, pp. 403–424, 1997.
- [35] M. Bishop, *Computer Security: Art and Science*. Reading, MA: Addison-Wesley Professional, 2002.
- [36] R. Chabukswar, Y. Mo, and B. Sinopoli, "Detecting integrity attacks on SCADA systems," in *Proc. 18th IFAC World Congr.*, Milano, Italy, Aug.–Sept. 2011, pp. 11239–11244.
- [37] (2011, Apr.). Risk management fundamentals. U.S. Dept. Homeland Security. [Online]. Available: www.dhs.gov/xlibrary/assets/rma-risk-management-fundamentals.pdf
- [38] (2012, Sept.). NIST special publication 800-30: Guide for conducting risk assessments. U.S. Dept. Commerce. Natl. Inst. Standards Technol., Revision 1. Tech. Rep. NIST SP 800-30 Rev. 1 [Online]. Available: csrc.nist.gov/publications/nistpubs/800-30-rev1/sp800_30_r1.pdf
- [39] S. Kaplan and B. J. Garrick, "On the quantitative definition of risk," *Risk Anal.*, vol. 1, no. 1, pp. 11–27, 1981.
- [40] T. Somestad, M. Ekstedt, and H. Holm, "The cyber security modeling language: A tool for assessing the vulnerability of enterprise system architectures," *IEEE Syst. J.*, vol. 7, no. 3, pp. 363–373, 2013.
- [41] I. Garitano, R. Uribeetxeberria, and U. Zurutuza, "A review of SCADA anomaly detection systems," in *Proc. 6th Int. Conf. SOCO Soft Computing Models Industrial Environmental Applications* (Advances in Intelligent and Soft Computing), E. Corchado, V. Šnášel, J. Sedano, A. E. Hassanien, J. L. Calvo, and D. Ślęzak, Eds. Berlin, Heidelberg, Germany: Springer, 2011, vol. 87, pp. 357–366.
- [42] A. Abur and A. Exposito, *Power System State Estimation: Theory and Implementation*. New York: Marcel-Dekker, 2004.
- [43] A. Monticelli, *State Estimation in Electric Power Systems: A Generalized Approach*. Boston, MA: Kluwer Academic, 1999.
- [44] Y. Liu, M. K. Reiter, and P. Ning, "False data injection attacks against state estimation in electric power grids," in *Proc. 16th ACM Conf. Computer Communications Security*, Chicago, IL, Nov. 2009, pp. 21–32.
- [45] G. Dán and H. Sandberg, "Stealth attacks and protection schemes for state estimators in power systems," in *Proc. IEEE Smart Grid Communications*, Oct. 2010, pp. 214–219.
- [46] R. Bobba, K. M. Rogers, Q. Wang, H. Khurana, K. Nahrstedt, and T. Overbye, "Detecting false data injection attacks on DC state estimation," in *Proc. Preprints 1st Workshop Secure Control Systems*, Stockholm, Sweden, Apr. 2010.
- [47] O. Kosut, L. Jia, R. Thomas, and L. Tong, "Malicious data attacks on smart grid state estimation: Attack strategies and countermeasures," in *Proc. 1st IEEE Int. Conf. Smart Grid Communications*, Gaithersburg, MD, Oct. 2010, pp. 220–225.
- [48] A. Giani, E. Bitar, M. McQueen, P. Khargonekar, K. Poola, and M. Garcia, "Smart grid data integrity attacks: Characterizations and countermeasures," in *Proc. IEEE Smart Grid Communications*, Oct. 2011, pp. 232–237.
- [49] T. T. Kim and H. V. Poor, "Strategic protection against data injection attacks on power grids," *IEEE Trans. Smart Grid*, vol. 2, pp. 326–333, June 2011.
- [50] A. Tillmann and M. Pfetsch. (2012). The computational complexity of the restricted isometry property, the nullspace property, and related concepts in compressed sensing. [Online]. Available: <http://arxiv.org/abs/1205.2081>
- [51] S. McCormick, "A combinatorial approach to some sparse matrix problems," Ph.D. dissertation, Stanford Univ., Stanford, CA, 1983.
- [52] S. Jökar and M. E. Pfetsch, "Exact and approximate sparse solutions of underdetermined linear equations," *SIAM J. Sci. Comput.*, vol. 31, no. 1, pp. 23–44, Oct. 2008.
- [53] J. Tsitsiklis and D. Bertsimas, *Introduction to Linear Optimization*. Belmont, MA: Athena Scientific, 1997.
- [54] A. Schrijver, *Theory of Linear and Integer Programming*. New York: Wiley, 1986.
- [55] A. M. Bruckstein, D. L. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *SIAM Rev.*, vol. 51, no. 1, pp. 34–81, 2009.
- [56] E. Candès and T. Tao, "Decoding by linear programming," *IEEE Trans. Inform. Theory*, vol. 51, no. 12, pp. 4203–4215, Dec. 2005.
- [57] B. Recht, M. Fazel, and P. A. Parrilo, "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization," *SIAM Rev.*, vol. 52, no. 3, pp. 471–501, 2010.
- [58] U.S.–Canada PSOTF. (2014, Apr.). Final report on the August 14th blackout in the United States and Canada. [Online]. Available: energy.gov/sites/prod/files/oeprod/DocumentsandMedia/BlackoutFinal-Web.pdf
- [59] J. Hendrickx, K. H. Johansson, R. Jungers, H. Sandberg, and K. C. Sou, "Efficient computations of a security index for false data attacks in power networks," *IEEE Trans. Automat. Contr.*, to be published, doi: 10.1109/TAC.2014.2351625.
- [60] G. R. Krumpholz, K. Clements, and P. Davis, "Power system observability: A practical algorithm using network topology," *IEEE Trans. Power Apparatus Syst.*, vol. PAS-99, no. 4, pp. 1534–1542, 1980.
- [61] K. C. Sou, H. Sandberg, and K. Johansson, "Computing critical k-tuples in power networks," *IEEE Trans. Power Syst.*, vol. 27, no. 3, pp. 1511–1520, 2012.
- [62] M. Stoer and F. Wagner, "A simple min-cut algorithm," *J. ACM*, vol. 44, pp. 585–591, July 1997.
- [63] K. C. Sou, H. Sandberg, and K. H. Johansson, "Data attack isolation in power networks using secure voltage magnitude measurements," *IEEE Trans. Smart Grid*, vol. 5, no. 1, pp. 14–28, Jan. 2004.
- [64] R. Christie. (1993). Power system test case archive. Univ. Washington. Seattle, WA. Tech. Rep. [Online]. Available: http://www.ee.washington.edu/research/pstca/pf118/pg_tca118bus.htm
- [65] R. T. Marler and J. S. Arora, "Survey of multi-objective optimization methods for engineering," *Structural Multidisciplinary Optim.*, vol. 26, no. 6, pp. 369–395, Apr. 2004.
- [66] K. Johansson, "The quadruple-tank process: A multivariable laboratory process with an adjustable zero," *IEEE Trans. Control Syst. Technol.*, vol. 8, no. 3, pp. 456–465, May 2000.
- [67] (2012, Apr.). Supplemental material: Experimental setup and attack scenarios. [Online]. Available: <http://urn.kb.se/resolve?urn=urn:nbn:se:kt:h:diva-96745>
- [68] P. Kundur, *Power System Stability and Control*. New York: McGraw-Hill Professional, 1994.
- [69] CPLEX. (2014, 17 Nov.). [Online]. Available: <http://www-01.ibm.com/software/commerce/optimization/cplex-optimizer/>
- [70] Gurobi Optimization, Inc. (2013). Gurobi optimizer reference manual. [Online]. Available: www.gurobi.com
- [71] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, *Network Flows: Theory, Algorithms, and Applications*. Englewood Cliffs, NJ: Prentice Hall, 1993.

