

A Divide-and-Conquer Approach to Distributed Attack Identification

Fabio Pasqualetti^a, Florian Dörfler^b, Francesco Bullo^c

^aDepartment of Mechanical Engineering, University of California, Riverside

^bDepartment of Electrical Engineering, University of California, Los Angeles

^cDepartment of Mechanical Engineering, University of California, Santa Barbara

Abstract

Identifying attacks is key to ensure security in cyber-physical systems. In this note we remark upon the computational complexity of the attack identification problem by showing how conventional approximation techniques may fail to identify attacks. Then, we propose decentralized and distributed monitors for attack identification with performance guarantees and low computational complexity. The proposed monitors rely on the geometric framework proposed in [1], yet require only local knowledge of the system dynamics and parameters. We exploit a *divide-and-conquer* approach, where first the system is partitioned into disjoint regions, then corrupted regions are identified via distributed computation, and finally corrupted components are isolated within regions.

Keywords: Security, cyber-physical system, attack identification, distributed control and estimation, network system.

1. Introduction

Cyber-physical systems are the core of many technological domains, including health care and biomedicine, telecommunications, and energy management. Due to their importance, cyber-physical systems are not only prone to sensor and actuator failures as legacy control systems, but also to intentional attacks against control and communications modules. Attacks can have major consequences, ranging from significant economic losses to instabilities and services disruption [2, 3, 4].

Detection and identification of attacks is key to design effective security mechanisms. Fundamental limitations in the detectability and identifiability of attacks have recently been characterized for different system dynamics, attack models, and monitoring systems. For instance, in [1, 5, 6, 7, 8] it is shown how attackers with access to sufficiently many system resources can always avoid detection and identification, as well as attackers with more limited resources and full knowledge of the system dynamics and state. Conversely, if the monitoring resources and information outbalance the attack capabilities, the attack locations and strategy can be promptly reconstructed. Moreover, while detecting attacks is computationally *easy* in both centralized and distributed settings [1, 9], identifying the attack location and strategy is computationally *hard* [1].

Despite its importance, few solutions have been proposed for the identification of attacks. A complete, yet computationally intensive, solution to the attack identification problem is described in [1] by using unknown-input observers and geometric

control techniques [10]. Convex relaxation techniques are employed in [11] to derive an efficient (yet incomplete and without guarantees) identification algorithm for the case of attacks against the system measurements. Finally, [12] shows that certain instances of the identification problem can in fact be solved efficiently. In this work we derive decentralized and distributed identification monitors with performance guarantees.

The main contributions of this note are as follows. First, we remark on the complexity of the attack identification problem and show how common convex relaxation techniques may fail to identify attacks (Section 2). Our examples highlight that, in large-scale systems, different output and state attacks may achieve the same cost in relaxed optimization problems, thereby impeding their correct identification. The inherent computational complexity and shortcoming of relaxation methods motivate our second contribution: we present a fully decentralized and low-complexity identification method and characterize its performance (Section 3.2). Our decentralized method relies on geographically distributed control centers, which have local knowledge of the system parameters. We show that the performance of our decentralized identification method depends only on the system structure and parameters, and not on the strategy of attack. Hence, our decentralized method also provides guidelines for the design of secure cyber-physical systems. Third, we propose a distributed identification method based on the *divide-and-conquer* principle (Section 3.3). Analogously to our decentralized method, our distributed algorithm requires only local model information and communication, and it achieves guaranteed identification of a class of attacks. Our distributed method overcomes the performance of its decentralized counterpart, at the expense of a more involved algorithmic structure. Finally, as a minor contribution, we present a state estimation algorithm for descriptor systems with unknown inputs (Appendix A).

Email addresses: fabiopas@engr.ucr.edu (Fabio Pasqualetti),
dorfler@seas.ucla.edu (Florian Dörfler),
bullo@engineering.ucsb.edu (Francesco Bullo)

URL: <http://www.fabiopas.it> (Fabio Pasqualetti),
<http://www.seas.ucla.edu/~dorfler/> (Florian Dörfler),
<http://motion.mee.ucsb.edu/> (Francesco Bullo)

2. The centralized attack identification problem

In this section we present our setup for the attack identification problem, and we recall some results and fundamental limitations of centralized identification methods.

2.1. Centralized setup and notation

We represent a cyber-physical system under attack with the continuous-time, linear, and time-invariant descriptor system ¹

$$\begin{aligned} E\dot{x}(t) &= Ax(t) + Bu(t), \\ y(t) &= Cx(t) + Du(t), \end{aligned} \quad (1)$$

where $x : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^n$, $y : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^p$, $E \in \mathbb{R}^{n \times n}$, $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$, and $D \in \mathbb{R}^{p \times m}$. The matrix E is possibly singular, and the inputs $Bu : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^n$ and $Du : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^p$ are unknown signals describing disturbances affecting the plant. Besides reflecting the genuine failure of systems components, these unknown inputs model the effect of attacks against cyber and physical components. We assume that each state and output variable can be independently compromised. Accordingly, we partition the input matrices and attack signals as $B = [I \ 0]$, $D = [0 \ I]$, and $u = [u_x^\top, u_y^\top]^\top$, where u_x and u_y are referred to as *state attack* and *output attack*, respectively. As shown in [1], many interesting cyber-physical systems and attacks can be modeled by the descriptor system (1) subject to unknown inputs.

The attack signal depends upon the attack strategy. In particular, if the *attack set* (or attacked variables) is $K \subseteq \{1, \dots, n+p\}$, with $|K| = k$, then only the entries of u indexed by K are nonzero over time, that is, for each $i \in K$, there exists a time t such that $u_i(t) \neq 0$, and $u_j(t) = 0$ for all $j \notin K$ and at all times. To underline this sparsity relation, we sometimes use u_K to denote the attack signal, that is the subvector of u indexed by K . Analogously, the pair (B_K, D_K) , where B_K and D_K are the submatrices of B and D with columns in K , are referred to as the *attack signature*. Hence, $Bu = B_K u_K$, and $Du = D_K u_K$.

We make the following assumptions on system (1):

- (A1) the pair (E, A) is regular, that is, the determinant $|sE - A|$ is nonzero for almost all values $s \in \mathbb{C}$;
- (A2) the initial condition $x(0) \in \mathbb{R}^n$ is consistent, that is, $(Ax(0) + Bu(0)) \in \text{Im}(E)$; and
- (A3) the input u is smooth.

Assumption (A1) ensures the existence of a unique solution $x(t)$ to (1). Assumptions (A2) and (A3) guarantee smoothness of the state trajectory and the measurements; see [13, Lemma 2.5]. If assumptions (A2) and (A3) are dropped, then there are inconsistent initial conditions and impulsive inputs by which a powerful attacker can avoid detection [1]. Finally, we assume that the cardinality k of the attack set, or an upper bound, is known.

¹The results stated in this paper for continuous-time descriptor systems hold also for discrete-time descriptor systems and nonsingular systems. Moreover, due to linearity of (1), known inputs do not affect our results and are not included in the model.

2.2. Identifiability of cyber-physical attacks

Informally, an attack K is unidentifiable if it cannot be distinguished (from knowledge of the measurements and the system parameters) from another attack R corrupting equally many or fewer variables. Here, we confine ourselves to comparing the attack set K with other attack sets R with $|R| \leq |K|$ since sufficiently large attack sets can always be designed to be unidentifiable, for instance, by corrupting sufficiently many sensors.

More formally, let $y(x_0, u, t)$ be the output sequence generated from the initial state x_0 under the attack signal u . We adopt the following definition of identifiability of attacks [1]:

Definition 1. (Identifiability of cyber-physical attacks) For the descriptor system (1) with initial state x_0 , the attack $(B_K u_K, D_K u_K)$ is unidentifiable if and only if $y(x_0, u_K, t) = y(x_1, u_R, t)$ for some initial state $x_1 \in \mathbb{R}^n$, for some attack $(B_R u_R, D_R u_R)$ with $|R| \leq |K|$ and $R \neq K$, and for all $t \in \mathbb{R}_{\geq 0}$.

In [1, Theorem 3.4], we provided the following equivalent system-theoretic characterization of identifiability:

Theorem 2.1. (Algebraic test for identifiability of cyber-physical attacks) For the descriptor system (1) and an attack set K , the following statements are equivalent:

- (i) the attack set K is unidentifiable; and
- (ii) there is an attack set R , with $|R| \leq |K|$ and $R \neq K$, $s \in \mathbb{C}$, $g_K \in \mathbb{C}^{|K|}$, $g_R \in \mathbb{C}^{|R|}$, and $x \in \mathbb{C}^n$, with $x \neq 0$, such that

$$\begin{aligned} (sE - A)x - \begin{bmatrix} B_K & B_R \end{bmatrix} \begin{bmatrix} g_K \\ g_R \end{bmatrix} &= 0, \\ Cx + \begin{bmatrix} D_K & D_R \end{bmatrix} \begin{bmatrix} g_K \\ g_R \end{bmatrix} &= 0. \end{aligned} \quad (2)$$

Condition (2) shows that existence of unidentifiable attack sets of cardinality k is equivalent to the existence of *invariant zeros* for the system $(E, A, B_{\bar{K}}, C, D_{\bar{K}})$ with $|\bar{K}| \leq 2k$. We refer to [10, 14] for a review invariant zeros of descriptor systems.

2.3. Centralized identification: complexity and pitfalls

The attack identification problem is concerned with identifying the attack set K from measurements y and knowledge of the system parameters (E, A, C) . The identification problem can be reformulated as the following cardinality minimization problem [1, Lemma 4.4]: given a descriptor system with matrices $E, A \in \mathbb{R}^{n \times n}$, measurement matrix $C \in \mathbb{R}^{p \times n}$, and measurement signal $y : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^p$, find the minimum cardinality input signals $v_x : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^n$ and $v_y : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^p$ and an initial condition $\xi_0 \in \mathbb{R}^n$ that explain the measurements y , that is,

$$\begin{aligned} \min_{v_x, v_y, \xi_0} \quad & \|v_x\|_{\mathcal{L}_0} + \|v_y\|_{\mathcal{L}_0} \\ \text{subject to} \quad & E\dot{\xi}(t) = A\xi(t) + v_x(t), \\ & y(t) = C\xi(t) + v_y(t), \\ & \xi(0) = \xi_0 \in \mathbb{R}^n. \end{aligned} \quad (3)$$

Here we use the shorthands $\text{supp}(x) = \{i \in \{1, \dots, n\} : x_i \neq 0\}$ for the number of non-zero entries of a vector $x \in \mathbb{R}^n$ and $\|v\|_{\mathcal{L}_0} = |\cup_{t \in \mathbb{R}_{\geq 0}} \text{supp}(v(t))|$ for a vector-valued signal $v : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^n$.

The optimization problem (3) is generally combinatorial and belongs to the class of *NP-hard* problems [1, Corollary 4.5].

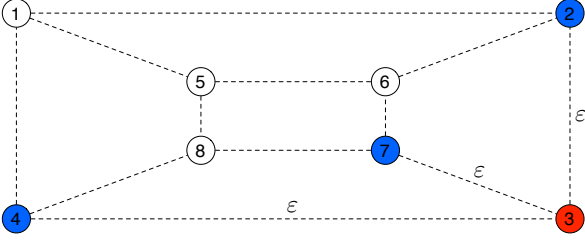


Figure 1: A regular consensus system (A, B, C) , where the state variable 3 is corrupted by the attacker, and the state variables 2, 4, and 7 are directly measured. Due to the sparsity pattern of (A, B, C) any attack of cardinality one is generically detectable and identifiable; see [1, 16] for further details.

Lemma 2.2. (Complexity of the attack identification problem) Consider the descriptor system (1) with identifiable attack set K . The attack identification problem given the system matrices E, A, C , and the measurements y is NP-hard.

As a consequence of this inherent complexity, the identification of the attack set K requires a combinatorial procedure, since, a priori, K is one of the $\binom{n+p}{|K|}$ possible attack sets. In [1, Section 4.D], the authors provided a solution based on the implementation of $\binom{n+p}{|K|}$ residual filters [10] each determining whether a predefined set coincides with the attack set. The solution in [1] is *complete*, but does not scale to large attack sets.

In the case of discrete-time and regular systems subject to output attacks, the attack identification problem can be solved efficiently if the monitoring system has access to a substantial amount of resources. The particular assumption is that the pair (A, C) remains observable after removing any set of $2|K|$ rows of C (that is, any set of $2|K|$ sensors) [12, Propositions 3.2 and 3.3]. If this strong observability assumption is not met, or in case of state attacks on (regular or singular) systems, the problem remains computationally hard. In this case, a natural approach is to apply convex relaxation approaches to the optimization problem (3). Cardinality minimization problems of the form $\min_{v \in \mathbb{R}^n} \text{supp}(y - Av)$ can often be efficiently solved via the ℓ_1 regularization $\min_{v \in \mathbb{R}^n} \|y - Av\|_{\ell_1}$ [15]. This procedure can be adapted to problem (3) after converting it into an algebraic problem, for instance by taking subsequent derivatives of the output y , or by discretizing the continuous-time system (1) and recording several measurements. As shown in [11], for discrete-time systems the ℓ_1 regularization performs reasonably well in the presence of output attacks. However, in the presence of state attacks such an ℓ_1 relaxation may perform poorly. In the following, we develop an intuition explaining why this approach may fail, particularly in large-scale systems.

Example 1. (Ineffectiveness of regularization methods for sufficiently distant attacks) Consider a consensus system with underlying network graph (reflecting the sparsity pattern of A) illustrated in Fig. 1. In our model (1), the system matrices are taken as $E = I$ and, for $0 < \varepsilon \ll 1$, A is the negative Laplacian

$$A = \begin{bmatrix} -0.8 & 0.1 & 0 & 0.2 & 0.5 & 0 & 0 & 0 \\ 0.1 & -0.4-\varepsilon & \varepsilon & 0 & 0 & 0.3 & 0 & 0 \\ 0 & 3\varepsilon & -9\varepsilon & 0 & 0 & 0 & 6\varepsilon & 0 \\ 0.1 & 0 & \varepsilon & -0.5-\varepsilon & 0 & 0 & 0 & 0.4 \\ 0.1 & 0 & 0 & 0 & -0.6 & 0.2 & 0 & 0.3 \\ 0 & 0.4 & 0 & 0 & 0.1 & -0.6 & 0.1 & 0 \\ 0 & 0 & 3\varepsilon & 0 & 0 & 0.4 & -0.6-3\varepsilon & 0.2 \\ 0 & 0 & 0 & 0.3 & 0.2 & 0 & 0.2 & -0.7 \end{bmatrix}.$$

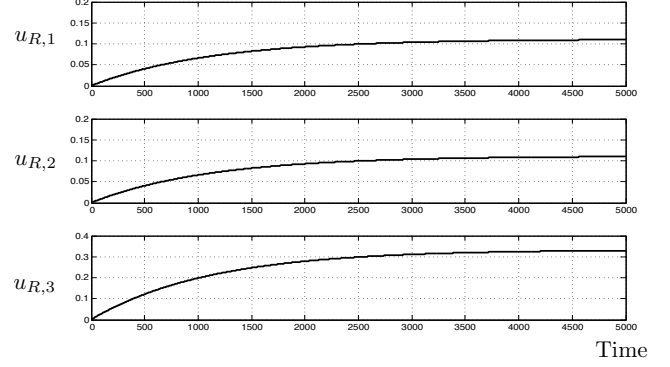


Figure 2: In Fig. 2 we plot the attack mode u_R for the attack set $R = \{2, 4, 7\}$ to generate the same output as the attack set $K = \{3\}$ with attack mode $u_K = 1$. Although $|R| > |K|$, we have that $|u_{R,i}(t)| < |u_K(t)|/3$ for $i \in \{1, 2, 3\}$.

Let the measurement matrix C and the attack signature B_K be

$$C = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}, \quad B_K^\top = [0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0],$$

and define the transfer matrix $G_K(s) = C(sI - A)^{-1}B_K$. It can be verified that the state attack $K = \{3\}$ is identifiable.

Consider also the state attack $R = \{2, 4, 7\}$ with signature

$$B_R^\top = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} = C,$$

and define the transfer matrix $G_R(s) = C(sI - A)^{-1}B_R$. Let $U_K(s)$ and $U_R(s)$ be the Laplace transforms of $u_K(t)$ and $u_R(t)$, respectively. Notice that $G_R(s)$ is right-invertible [10]. Thus,

$$Y(s) = G_K(s)U_K(s) = G_R(s)(G_R^{-1}(s)G_K(s)U_K(s)).$$

In other words, the measurements $Y(s)$ generated by the attack $U_K(s)$ can equivalently be generated by the attack

$$U_R(s) = G_R^{-1}(s)G_K(s)U_K(s).$$

Notice that $3 = \|u_R\|_{\mathcal{L}_0} > \|u_K\|_{\mathcal{L}_0} = 1$, that is, the attack set K achieves a lower cost than R in the optimization problem (3).

Consider now the numerical realization $\varepsilon = 0.0001$, $x(0) = 0$, and $u_K(t) = 1$ for all $t \in \mathbb{R}_{\geq 0}$. The corresponding attack signal $u_R = [u_{R,1} \ u_{R,2} \ u_{R,3}]$ is shown in Fig. 2. Observe that

$$\|u_K(t)\|_{\ell_p} > \|u_R(t)\|_{\ell_p}$$

holds point-wise in time for all integers $p \geq 1$. We also have

$$\|u_K(t)\|_{\mathcal{L}_q/\ell_p} > \|u_R(t)\|_{\mathcal{L}_p/\ell_q}$$

for any integers $p, q \geq 1$ and with the \mathcal{L}_q/ℓ_p signal norm

$$\|u_K(t)\|_{\mathcal{L}_q/\ell_p} = \left(\int_0^\infty \|u_K\|_p^q d\tau \right)^{1/q}.$$

Hence, the attack set R achieves a lower cost than K for any version of the optimization problem (3) penalizing a ℓ_p cost point-wise in time or a \mathcal{L}_q/ℓ_p cost over a time interval. On the other hand, we have $\|u_R\|_{\mathcal{L}_0} > \|u_K\|_{\mathcal{L}_0}$. We conclude that, in

general, the identification problem cannot be solved by a point-wise ℓ_p or \mathcal{L}_q/ℓ_p regularization for any $p, q \geq 1$. Finally, we remark that for any choice of network parameters, a value of ε can be found such that a point-wise ℓ_p or a \mathcal{L}_q/ℓ_p regularization procedure fails at identifying the attack set. \square

We emphasize that Example 1 is not of pathological nature, but large-scale stable systems often exhibit this behavior independently of the system parameters for attacks which are “sufficiently distant” from the sensors. This can be easily seen in regular discrete-time systems, where a state attack with attack set K affects the output via the matrix $CA^{r-1}B_K$, where r is the relative degree of (A, B_K, C) . Notice that if A is Schur stable then $\lim_{k \rightarrow \infty} A^k = 0$, and $CA^{r-1}B_K$ converges to the zero matrix for increasing relative degree. In this case, an attack closer to the sensors may achieve a lower \mathcal{L}_q/ℓ_p cost than an attack far from sensors independently of the cardinality of the attack set. In short, the ε -connections in Example 1 can be thought of as the effect of a large relative degree in large-scale systems.

We conclude that centralized attack identification procedures are generally not tractable due to their inherent combinatorial complexity. Special cases that are tractable using efficient and provably correct procedures require restrictive observability assumptions. Finally, naive convex relaxation approaches often fail in large-scale systems with sufficiently distant state attacks.

3. The distributed attack identification problem

The obstacles and pitfalls in the centralized attack identification problem motivate our study of *divide-and-conquer* methods. In this section, we design distributed attack identification algorithms with performance guarantees, requiring low computational cost and local knowledge of the system parameters.

3.1. Distributed setup and notation

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be the undirected graph associated with the pair (E, A) , where the vertex set $\mathcal{V} = \{1, \dots, n\}$ corresponds to the system states, and the set of edges $\mathcal{E} = \{(i, j) : e_{ij} \neq 0 \text{ or } a_{ij} \neq 0 \text{ or } e_{ji} \neq 0 \text{ or } a_{ji} \neq 0\}$ is induced by the sparsity pattern of E and A ; see also [1, Section IV]. Assume that \mathcal{V} is partitioned into N disjoint subsets as $\mathcal{V} = \mathcal{V}_1 \cup \dots \cup \mathcal{V}_N$, with $|\mathcal{V}_i| = n_i$, and let $\mathcal{G}^i = (\mathcal{V}_i, \mathcal{E}_i)$ be the i -th subgraph of \mathcal{G} with vertices \mathcal{V}_i and edges $\mathcal{E}_i = \mathcal{E} \cap (\mathcal{V}_i \times \mathcal{V}_i)$. According to this partition, and possibly after relabeling the states, the system matrix A in (1) can be written as

$$A = \begin{bmatrix} A_1 & \cdots & A_{1N} \\ \vdots & \vdots & \vdots \\ A_{N1} & \cdots & A_N \end{bmatrix} = A_D + A_C,$$

where $A_i \in \mathbb{R}^{n_i \times n_i}$, $A_{ij} \in \mathbb{R}^{n_i \times n_j}$, A_D is block-diagonal, and $A_C = A - A_D$. Notice that, if $A_D = \text{blkdiag}(A_1, \dots, A_N)$, then A_D represents the isolated subsystems and A_C describes the interconnection structure among the subsystems. Additionally, if the original system is sparse, then several blocks in A_C vanish.

We make the following assumptions on the subsystem decomposition:

- (A4) the matrices E and C are block-diagonal, that is $E = \text{blkdiag}(E_1, \dots, E_N)$ and $C = \text{blkdiag}(C_1, \dots, C_N)$, where $E_i \in \mathbb{R}^{n_i \times n_i}$ and $C_i \in \mathbb{R}^{p_i \times n_i}$,
- (A5) each pair (E_i, A_i) is regular, and each triple (E_i, A_i, C_i) is observable.

Let $\mathcal{N}_i = \{j \in \{1, \dots, N\} \setminus \{i\} : \|A_{ij}\| \neq 0 \text{ or } \|A_{ji}\| \neq 0\}$ be the neighbors of subsystem i , and let \mathcal{N}_i^k be the set of neighbors at distance k from i , with subsystem i excluded. Each subsystem \mathcal{G}^i has a *control center* with the following capabilities:

- (A6) the i -th control center knows the matrices E_i, A_i, C_i , as well as the neighboring matrices A_{ij} , $j \in \mathcal{N}_i$; and
- (A7) the i -th control center can transmit an estimate of its state to the j -th control center if $j \in \mathcal{N}_i$.

Given the above structure, the descriptor system (1) can be written as the interconnection of N subsystems of the form

$$\begin{aligned} E_i \dot{x}_i(t) &= A_i x_i(t) + \sum_{j \in \mathcal{N}_i} A_{ij} x_j(t) + B_{K_i} u_{K_i}(t), \\ y_i(t) &= C_i x_i(t) + D_{K_i} u_{K_i}(t), \quad i \in \{1, \dots, N\}, \end{aligned} \quad (4)$$

Here $K_i = (K \cap \mathcal{V}_i) \cup K_i^p$ is the attack set in region \mathcal{G}^i , and K_i^p is the set of corrupted measurements in region \mathcal{G}^i . Clearly, if the inter-subsystem signals $A_{ij} x_j$ are known or directly measured, then the regional attack identification problem within each subsystem reduces to the centralized problem. In the following, we will not make this assumption since it is restrictive and implicitly precludes the case that the inter-subsystem signals $A_{ij} x_j$ themselves are corrupted by an attacker.

3.2. Fully decoupled attack identification

As a first low-complexity identification method we consider the fully decoupled case (no cooperation among control centers). In the spirit of fully decentralized state estimation [17], the neighboring states x_j affecting x_i are treated as unknown inputs f_i to the i -th subsystem, and equation (4) becomes

$$\begin{aligned} E_i \dot{x}_i(t) &= A_i x_i(t) + B_i^b f_i(t) + B_{K_i} u_{K_i}(t), \\ y_i(t) &= C_i x_i(t) + D_{K_i} u_{K_i}(t), \quad i \in \{1, \dots, N\}, \end{aligned} \quad (5)$$

where $B_i^b = [A_{i1} \cdots A_{i,i-1} A_{i,i+1} \cdots A_{iN}]$. We refer to (5) as to the *i -th decoupled system*, and we let $\mathcal{V}_i^b \subseteq \mathcal{V}_i$ be the set of *boundary nodes* of (5), that is, the nodes $j \in \mathcal{V}_i$ with $A_{jk} \neq 0$ for some $k \in \{1, \dots, n\} \setminus \mathcal{V}_i$. Due to partitioning, control centers perform attack identification only on local subsystems.

To explicitly identify the attack set K_i we construct a residual generator that is insensitive to the inputs $B_i^b f_i$ and $B_{K_i} u_{K_i}$. Following our work [1], we resort to a geometric control approach [10] and construct a residual filter of the form

$$\begin{aligned} E_i \dot{w}_i(t) &= (A_i + L_i C_i) w_i(t) - L_i y_i(t), \\ r_i(t) &= M w_i(t) - H y_i(t), \end{aligned} \quad (6)$$

where L_i is the injection matrix associated with the conditioned invariant subspace generated by $[B_i^b \ B_{K_i}]$ and such that $(E_i, A_i + L_i C_i)$ is Hurwitz. The matrices M and H in (6) are chosen so that the inputs $B_i^b f_i + B_{K_i} u_{K_i}$ do not affect the residual r_i . In

summary, r_i is not identically zero if and only if the inputs to (6) are linearly independent from $B_i^b f_i + B_{K_i} u_{K_i}$. We refer to [1] for a detailed construction and discussion of this type of filter.

To uniquely identify the attack set K_i affecting region \mathcal{G}^i , each control center needs to construct one residual filter of the form (6) for each attack set of cardinality $|K_i|$. Thus, compared to the centralized case requiring one residual filter for each attack set of size $|K| = \sum_{i=1}^N |K_i|$, the computational complexity of the attack identification problem is tremendously reduced. On the other hand, some fundamental limitations arise by naively treating the neighboring signals as unknown inputs.

Theorem 3.1. (Fully decoupled attack identification) *For the partitioned descriptor system (5) and an attack set K , the following statements are equivalent:*

- (i) *the attack set K in (5) is unidentifiable by the fully decoupled identification algorithm (6); and*
- (ii) *for some region $i \in \{1, \dots, N\}$ with $K_i \neq \emptyset$, there exists an attack set R_i , with $|R_i| \leq |K_i|$ and $R_i \neq K_i$, so that the system $(E, A, [B_i^b \ B_{K_i} \ B_{R_i}], C, [D_{K_i} \ D_{R_i}])$ has invariant zeros.*

Proof. Let $y_i(x_{i,0}, u_{K_i}, f_i, t)$ denote the output of system (5) at time t , with initial value $x_{i,0}$, attack input u_{K_i} , and boundary input f_i . Notice that the attack set K_i is undistinguishable from R_i if and only if $y(x_{i,0}, u_{K_i}, f_i, t) = y(x_{i,1}, u_{R_i}, h_i, t)$ at all times t , for some initial conditions $x_{i,0}$ and $x_{i,1}$, attack inputs u_{K_i} and u_{R_i} , and boundary inputs f_i and h_i . From Theorem 2.1 and due to linearity of the system, K_i is unidentifiable if and only if $y(x_{i,0} - x_{i,1}, u_{K_i} - u_{R_i}, f_i - h_i, t) = 0$ at all times. The claimed statement follows from the definition of invariant zeros [10] and the fact that the identification algorithm via the residual filter (6) is complete [1]. \square

By comparing Theorems 2.1 and 3.1 we conclude that the i -th control center cannot distinguish between an unknown input from a safe subsystem, an unknown input from a corrupted subsystem, and a boundary attack with the same input direction.

Corollary 3.2 (Limitation of decoupled algorithm). *The following statements hold for the partitioned descriptor system (5) with the fully decoupled identification algorithm (6):*

- (L1) *any (boundary) attack set $K_i \subseteq V_i^b$ is not identifiable by the i -th control center (in fact K_i is not detectable²), and*
- (L2) *any (external) attack set $K \setminus K_i$ is not identifiable by the i -th control center (in fact K_i is not detectable).*

3.3. Cooperative attack identification

In this section we improve upon the fully decoupled method presented in Subsection 3.2 and propose an identification method based on a *divide-and-conquer* procedure with cooperation. Our cooperative identification method is informally described as follows. First, control centers independently estimate the state of their local region subject to unknown inputs from the neighboring regions. Because of the presence of unknown inputs, the estimation computed by a control center is

correct modulo some *uncertainty subspace*. Control centers exchange their estimate and the corresponding uncertainty subspaces. Second, control centers check the compatibility between their estimate and those received from the neighboring regions. Third, if the received estimates are not compatible with local estimates, then the system is recognized under attack. Finally, control centers implement a local attack identification procedure by leveraging local system parameters and estimates, and estimates received from their neighbors.

We next detail our cooperative identification method.

(S1: local state estimation) Each control center estimates the state of its own region by means of an *unknown-input observer* for the i -th subsystem subject to the unknown input $B_i^b f_i$. We refer the reader to [10] for a detailed review of unknown-input observers, and to Appendix A for a constructive procedure.

Assume that the state x_i can be reconstructed modulo some subspace \mathcal{F}_i .³ Let $F_i = \text{Basis}(\mathcal{F}_i)$ be the uncertainty matrix, and partition the state accordingly as

$$x_i = \hat{x}_i + \tilde{x}_i, \quad (7)$$

where $\hat{x}_i(t) \perp \mathcal{F}_i$ is the portion of the state that can be estimated by the i -th control center in the presence of the unknown input $B_i^b f_i$, and $\tilde{x}_i(t) \in \mathcal{F}_i$. Let $z_i(t)$ be the *estimate* at time t of \hat{x}_i . Notice that, if the i -th region is not corrupted, then $z_i(t) = \hat{x}_i(t)$, whereas it may be $z_i(t) \neq \hat{x}_i(t)$ when $K_i \neq \emptyset$.

(S2: communication) Control centers transmit their estimate \hat{x}_i and uncertainty matrix F_i to every neighboring control center.

(S3: residual generation) Observe that

$$A_{ij}x_j = A_{ij}\hat{x}_j + A_{ij}\tilde{x}_j,$$

where \hat{x}_j and \tilde{x}_j are defined as in (7). After carrying out step (S1), since the matrices A_{ij} are known to the i -th control center due to Assumption (A6), only the inputs $A_{ij}\tilde{x}_j$ are unknown to the i -th control center, while the inputs $A_{ij}\hat{x}_j$ are known to the i -th center due to communication. Let

$$B_i^b F_i = [A_{i1}F_1 \ \cdots \ A_{i,i-1}F_{i-1} \ A_{i,i+1}F_{i+1} \ \cdots \ A_{iN}F_N],$$

and rewrite the signal $B_i^b \tilde{x}$ as $B_i^b \tilde{x} = B_i^b F_i f_i$, for some unknown signal f_i . Then the dynamics of the i -th subsystem read as

$$E_i \dot{\hat{x}}_i(t) = A_i x_i(t) + B_i^b \hat{x}_i(t) + B_i^b F_i f_i(t) + B_{K_i} u_{K_i}(t),$$

where \hat{x} is the vector of \hat{x}_i for all $i \in \{1, \dots, N\}$.

Next, we construct a residual generator akin to (6) that is insensitive to the input $B_i^b F_i f_i$ and makes use of the state estimates z transmitted to control center i by its neighbors:

$$\begin{aligned} E_i \dot{w}_i(t) &= (A_i + L_i C_i) w_i(t) - Ly(t) + B_i^b z(t), \\ r_i(t) &= M w_i(t) - Hy(t). \end{aligned} \quad (8)$$

Here L_i is the injection matrix associated with the conditioned invariant subspace generated by $B_i^b F_i$ and such that $(E_i, A_i +$

²An attack is detectable if it can be distinguished from the zero attack [1].

³For nonsingular systems without feedthrough matrix, \mathcal{F}_i is as small as the largest (A_i, B_i^b) -controlled invariant subspace contained in $\text{Ker}(C_i)$ [10].

Algorithm 1: Cooperative attack identification

Input : Matrices $E_i, A_i, A_{i,j}$ for $j \in \mathcal{N}_i$;
Require : Conditions (i), (ii), and (iii) in Theorem 3.3;
Output : Attack set K_i ;

- 1 Compute the uncertainty subspace $\mathcal{F}_i = \text{Im}(F_i)$;
- 2 Transmit F_i to control centers \mathcal{N}_i ;
- while True do**
- 3 Estimate state \hat{x}_i (state x modulo \mathcal{F}_i);
- 4 Transmit \hat{x}_i to \mathcal{N}_i , and receive \hat{x}_j from \mathcal{N}_i ;
- 5 Compute residual r_i as in (8);
- 6 Transmit r_i to \mathcal{N}_i , and receive r_j from \mathcal{N}_i ;
- 7 **if** $r_i \neq 0$ or $r_j \neq 0$ for all $j \in \mathcal{N}_i$ **then**
- 8 Identify K_i in local subsystem;
- 8 **return** K_i

$L_i C_i$) is Hurwitz. The matrices M and H in (8) are chosen so that the input $B_i^b F_i f_i$ does not affect the residual r_i [1].

(S4: cooperative attack identification) Neighboring control centers exchange the zero/nonzero status of the previously computed residuals, identify corrupted regions, and independently identify attacks in each attacked region. Our cooperative identification procedure for the i -th control center is summarized in Algorithm 1. We make the following technical assumptions:

- (A8) corrupted regions have one neighbor at distance 2, that is, $|\mathcal{N}_i^2| \geq 1$ for all regions i with $K_i \neq \emptyset$, and
- (A9) corrupted regions are separated by 3 non-corrupted regions, that is, $K_j = \emptyset$ for all $j \in \mathcal{N}_i^3$ and i with $K_i \neq \emptyset$.

Assumption (A8) requires a sufficiently large number of clusters, while assumption (A9) restricts the effectiveness of our procedure to geographically localized attacks. The next theorem characterizes the effectiveness of our cooperative identification procedure.

Theorem 3.3. (Cooperative attack identification) *For the partitioned system (4), the attack set K is identifiable by the cooperative identification algorithm if the following conditions hold:*

- (i) every system (E_i, A_i, B_i^b, C_i) has no invariant zeros, and
- (ii) every the system $(E_i, A_i, [B_i^b F_i B_{K_i} B_{R_i}], C_i, [D_{K_i} D_{R_i}])$ has no invariant zeros for every attack set R_i with $|R_i| \leq |K_i|$.

In Theorem 3.3, conditions (i) with assumptions (A8) and (A9) ensure *regional identifiability*, that is, the possibility to identify corrupted regions from local measurements and communication with neighboring regions. Condition (ii) ensures *local identifiability*, that is, attack identifiability within each corrupted region from local measurements and communication with neighboring regions. We defer the proof to Appendix B.

We conclude this section with the following observations. First, the cooperative identification procedure is implemented only on the corrupted regions (line 7 in Algorithm 1). Thus, the combinatorial complexity of our distributed identification procedure is $\sum_{i=1}^{\ell} \binom{n_i+p_i}{|K_i|}$, where ℓ is the number of corrupted regions. Hence, the distributed identification procedure greatly reduces the combinatorial complexity of the centralized procedure presented in [1] that requires the implementation of $\binom{n+p}{|K|}$

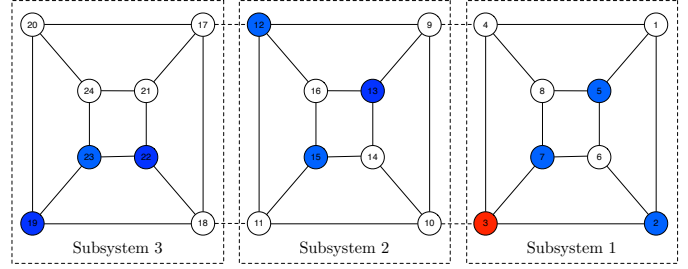


Figure 3: This figure shows a network composed of three subsystems. A control center is assigned to each subsystem. Each control center knows the dynamics of its local subsystem. The state of the blue nodes $\{2, 5, 7, 12, 13, 15\}$ is continuously measured by the corresponding control center, and the state of the red node $\{3\}$ is corrupted by an attacker. The decoupled identification procedure presented in Subsection 3.2 fails at detecting the attack. Instead, our cooperative identification procedure identifies the corrupted agent.

filters. Second, the conditions in Theorem 3.3 for cooperative identification improve upon those in Theorem 3.1 for fully decoupled identification; see Section 4 for an example. Third, our cooperative identification procedure is effective when attacks are localized in some regions, and regions under attack are sufficiently far from each other. Under these assumptions, our cooperative identification overcomes the limitations described in Example 1, because it does not rely on the magnitude of the measurements, and has provable performance guarantees.

4. Illustrative example

We now present an example showing that, contrary to the limitations of the naive fully decoupled approach (see Corollary 3.2), boundary attacks $K_i \subseteq V_i^b$ can be identified by our cooperative attack identification method.

Consider the sensor network in Fig. 3, where the state of the blue nodes $\{2, 5, 7, 12, 13, 15, 19, 22, 23\}$ is measured and the state of the red node $\{3\}$ is corrupted by an attacker. Assume that the network evolves according to nonsingular, linear, time-invariant dynamics. Assume further that the network has been partitioned into the three areas $\mathcal{V}_1 = \{1, \dots, 8\}$, $\mathcal{V}_2 = \{9, \dots, 16\}$, and $\mathcal{V}_3 = \{17, \dots, 24\}$. Since $\{3, 4\}$ are the boundary nodes for the first area, the attack set $K = 3$ is not identifiable (in fact it is not detectable) via the fully decoupled procedure in Section 3.2; see Corollary 3.2. It can be verified that the conditions in Theorem 3.3 are verified for generic network parameters [1, Section III.B], and that in fact the attack can be identified via our cooperative identification procedure. We conclude that our cooperative identification algorithm outperforms the decoupled identification algorithm in Section 3.2.

5. Conclusion

The problem of identifying attacks in cyber-physical systems requires a substantial computational effort. This paper shows that standard relaxation techniques may fail to identify state attacks in cyber-physical systems, and proposes two distributed algorithms with performance guarantees for attack identification by a set of geographically deployed control centers. The

algorithms require local measurements, local knowledge of the system, and communication with neighboring control centers. This paper provides initial results on the distributed attack identification problem, highlights its challenges and limitations, and foster the adoption of geometric control techniques for the solution of distributed control and estimation problems.

Appendix A. State estimation with unknown-input

In this section we present an algebraic technique to reconstruct the state of a descriptor system. Our method builds upon the results presented in [18]. Consider the descriptor model (1) written in the semi-explicit form (see [1, Section IV.C])

$$\begin{aligned} \dot{x}_1(t) &= A_{11}x_1(t) + A_{12}x_2(t) + B_1u(t), \\ 0 &= A_{21}x_1(t) + A_{22}x_2(t) + B_2u(t), \\ y(t) &= C_1x_1(t) + C_2x_2(t) + Du(t). \end{aligned} \quad (\text{A.1})$$

We aim at characterizing the largest subspace of the state space of (A.1) that can be reconstructed through the measurements y . Consider the associated nonsingular system

$$\begin{aligned} \dot{\tilde{x}}_1(t) &= A_{11}\tilde{x}_1(t) + B_1\tilde{u}(t) + A_{12}\tilde{x}_2(t), \\ \tilde{y}(t) &= \begin{bmatrix} \tilde{y}_1(t) \\ \tilde{y}_2(t) \end{bmatrix} = \begin{bmatrix} A_{21} \\ C_1 \end{bmatrix} \tilde{x}_1(t) + \begin{bmatrix} A_{22} & B_2 \\ C_2 & D \end{bmatrix} \begin{bmatrix} \tilde{x}_2(t) \\ \tilde{u}(t) \end{bmatrix}. \end{aligned} \quad (\text{A.2})$$

Recall from [10, Chapter 4] that the state of the system (A.2) can be reconstructed modulo its largest controlled invariant subspace \mathcal{V}_1^* contained in the null space of the output matrix.

Lemma Appendix A.1. (Reconstruction of the state $x_1(t)$) *Let \mathcal{V}_1^* be the largest controlled invariant subspace of the system (A.2). The state x_1 of the system (A.1) can be reconstructed only modulo \mathcal{V}_1^* through the measurements y .*

Proof. We start by showing that for every $x_1(0) \in \mathcal{V}_1^*$ there exist x_2 and u such that y is identically zero. Due to linearity of (A.1), we conclude that the projection of x_1 onto \mathcal{V}_1^* cannot be reconstructed. Notice that for every $\tilde{x}_1(0)$, \tilde{x}_2 , and \tilde{u} yielding $\tilde{y}_1 = 0$ at all times, the state trajectory $[\tilde{x}_1 \ \tilde{x}_2]$ is a solution to (A.1) with input $u = \tilde{u}$ and output $y = \tilde{y}_2$. Since for every $\tilde{x}_1(0) \in \mathcal{V}_1^*$ there exists \tilde{x}_2 and \tilde{u} such that \tilde{y} is identically zero, every state $x_1(0) \in \mathcal{V}_1^*$ cannot be reconstructed.

To conclude the proof, let $x_1(0)$ be orthogonal to \mathcal{V}_1^* , and let $x_1(t)$, $x_2(t)$, and $y(t)$ be the solution to (A.1) subject to the input $u(t)$. Notice that $\tilde{x}_1(t) = x_1(t)$, $\tilde{y}_1(t) = 0$, and $\tilde{y}_2(t) = y(t)$ is the solution to (A.2) with inputs $\tilde{x}_2(t) = x_2(t)$ and $\tilde{u}(t) = u(t)$. In other words, \tilde{x}_1 is a feasible state trajectory of (A.2) with inputs \tilde{x}_2 and \tilde{u} , and output \tilde{y} . By definition of \mathcal{V}_1^* , the state $\tilde{x}_1(0) \perp \mathcal{V}_1^*$ can be reconstructed from the measurements \tilde{y} [10]. \square

In the previous lemma we show that the state $x_1(t)$ of (A.1) can be reconstructed modulo \mathcal{V}_1^* . We now show that the state $x_2(t)$ can generally not be completely reconstructed.

Lemma Appendix A.2. (Reconstruction of the state $x_2(t)$) *Let $\mathcal{V}_1^* = \text{Im}(V_1)$ be the largest controlled invariant subspace of the system (A.2). The state x_2 of the system (A.1) can be reconstructed only modulo $\mathcal{V}_2^* = A_{22}^{-1} \text{Im}([A_{21} \ V_1 \ B_2])$.*

Proof. Let $x_1 = \bar{x}_1 + \hat{x}_1$, where $\bar{x}_1 \in \mathcal{V}_1^*$ and \hat{x}_1 is orthogonal to \mathcal{V}_1^* . From Lemma Appendix A.1, the signal \hat{x}_1 can be entirely reconstructed via y . Notice that

$$\begin{aligned} 0 &= A_{21}x_1(t) + A_{22}x_2(t) + B_2u(t), \\ &= A_{21}V_1v_1(t) + A_{21}\hat{x}_1(t) + A_{22}x_2(t) + B_2u(t), \end{aligned}$$

where V_1 is a basis of \mathcal{V}_1^* and $\bar{x}_1 = V_1v_1$. Let W be such that $\text{Ker}(W) = \text{Im}([A_{21}V_1 \ B_2])$. Then, $0 = WA_{21}\hat{x}_1(t) + WA_{22}x_2(t)$, and hence $x_2(t) = \bar{x}_2(t) + \hat{x}_2(t)$, where $\hat{x}_2(t) = (WA_{22})^\dagger WA_{21}\hat{x}_1(t)$, and $\bar{x}_2(t) \in \text{Ker}(WA_{22}) = A_{22}^{-1} \text{Im}([A_{21}V_1 \ B_2])$. The statement follows. \square

We remark that our characterization of \mathcal{V}_1^* and \mathcal{V}_2^* is equivalent to the definition of *weakly unobservable* subspace in [13], and of maximal *output-nulling* subspace in [19]. Hence, we proposed an optimal state estimator for our distributed attack identification procedure, and the matrix V_i in (**S1: local state estimation**) can be computed as in [13, 19]. Additionally, a reconstruction of x_1 modulo \mathcal{V}_1^* and x_2 modulo \mathcal{V}_2^* can be obtained through standard algebraic techniques [10]. Finally, the presented lemmas extend the results in [18] by characterizing the subspaces of the state space that can be reconstructed with an algebraic method by processing the measurements y and their derivatives.

Appendix B. Proof of Theorem 3.3

Before proving Theorem 3.3, we state some preliminary results. For the filter (8), define the error $e_i = w_i - x_i$ and note that

$$\begin{aligned} E_i \dot{e}_i(t) &= (A_i + L_i C_i) e_i(t) + B_i^b(z(t) - \hat{x}(t)) - B_{K_i} u_{K_i}(t) \\ &\quad - B_i^b F_i f_i(t), \\ r_i(t) &= M e_i(t), \end{aligned} \quad (\text{B.1})$$

where \hat{x} is the vector of all \hat{x}_i , and z_i is the estimate of \hat{x}_i by the j -th control center, and z is the vector of all z_i . We next show two fundamental properties of the residual r_i .

Lemma Appendix B.1. (Residual of isolated non-corrupted regions) *If $K_i = \emptyset$ and $K_j = \emptyset$ for all $j \in \mathcal{N}_i$, then r_i is identically zero.*

Proof. Consider a region j with $K_j = \emptyset$. Notice that the state estimation z_j satisfies $z_j = \hat{x}_j$. Because $K_i = \emptyset$, from (B.1) we conclude that the residual r_i is driven only by the input $B_i^b F_i f_i$. Since the residual generator (8) is constructed so that r_i is insensitive to the signature $(B_i^b F_i, 0)$, the statement follows. \square

Lemma Appendix B.2. (Residual of isolated corrupted regions) *For the partitioned system (4), let $K_i \neq \emptyset$. If*

- (i) $K_j = \emptyset$ for all $j \in \mathcal{N}_i^2$,
- (ii) every system (E_j, A_j, B_j^b, C_j) has no invariant zeros for all $j \in \mathcal{N}_i$, and
- (iii) the system $(E_i, A_i, [B_i^b F_i \ B_{K_i} \ B_{R_i}], C_i, [D_{K_i} \ D_{R_i}])$ has no invariant zeros for every attack set $R_i \neq K_i$ with $|R_i| \leq |K_i|$,

then

- (i) $r_i(t) \neq 0$ at some time $t \in \mathbb{R}_{\geq 0}$, and
- (ii) either $r_j(t) = 0$ for all $j \in \mathcal{N}_i$ at all times $t \in \mathbb{R}_{\geq 0}$, or $r_j(t) \neq 0$ for all $j \in \mathcal{N}_i$ at some times $t \in \mathbb{R}_{\geq 0}$.

Proof. The estimation computed by a control center is correct if its area is not under attack. In other words, since $K_j = \emptyset$ for all $j \in \mathcal{N}_i$, it follows $B_i^b \hat{x} = B_i^b z$ in (B.1). Because $(E_i, A_i, [B_i^b F_i B_{K_i}], C_i, [D_{K_i} D_{R_i}])$ has no invariant zeros, the attack set K_i is locally identifiable via local measurements and transmitted estimates, and statement (i) follows; see also [1].

To show the second statement, observe that only two cases are possible: either $\hat{x}_i = z_i$, or $\hat{x}_i \neq z_i$, where \hat{x}_i is defined in (7), and z_i is the estimate of \hat{x}_i computed by the i -th control center. For instance, if $\text{Im}(B_{K_i}) \subseteq \text{Im}(B_i^b)$, that is, the attack set K_i lies on the boundary of the i -th area, then $\hat{x}_i(t) = z_i(t)$.

In the first case, $\hat{x}_i = z_i$, all residuals r_j , $j \in \mathcal{N}_i$, are identically zero. In fact, since $K_\ell = \emptyset$ for all $\ell \in \mathcal{N}_i^2$, it follows that $\hat{x}_p = z_p$ for all $p \in \mathcal{N}_j$ and $j \in \mathcal{N}_i$, so that the j -th residual filter (B.1) evolves as an autonomous system, and r_j is identically zero.

Consider now the second case: $\hat{x}_i \neq z_i$. Notice that $B_j^b F_j f_j + B_j^b (\hat{x} - z) \in \text{Im}(B_j^b)$. Since (E_j, A_j, B_j^b, C_j) has no invariant zeros, every residual r_j , with $j \in \mathcal{N}_i$, cannot be identically zero. \square

We are now ready to prove Theorem 3.3.

Proof. Consider the i -th region, and let $K_i \neq \emptyset$. Due to conditions (i) and (ii) in Theorem 3.3, assumptions (A8) and (A9), and Lemma Appendix B.2 we conclude that:

- (C1) the residual r_i is not identically zero, and
- (C2) either $r_j(t) = 0$ for all $j \in \mathcal{N}_i$ at all times $t \in \mathbb{R}_{\geq 0}$, or $r_j(t) \neq 0$ for all $j \in \mathcal{N}_i$ at some times $t \in \mathbb{R}_{\geq 0}$.

Consider the region $p \in \mathcal{N}_i^2 \setminus \mathcal{N}_i$. Due to assumption (A9) and Lemma Appendix B.1 we conclude that:

- (C3) r_p is identically zero.

Consider the region $j \in \mathcal{N}_i$. Due assumption (A8) and the facts (C1) and (C3), we conclude that:

- (C4) there exists $j_1, j_2 \in \mathcal{N}_j$ such that r_{j_1} is identically zero, while r_{j_2} is not identically zero (take $j_1 = p$ and $j_2 = i$).

Corrupted regions are uniquely identified as the regions satisfying (C1) and (C2). See Figure B.4 for an example. Finally, due to condition (ii) in Theorem 3.3 each set K_i is locally identifiable (see also Theorem 2.1), and the statement follows. \square

References

- [1] F. Pasqualetti, F. Dörfler, F. Bullo, Attack detection and identification in cyber-physical systems, *IEEE Transactions on Automatic Control* 58 (11) (2013) 2715–2729.
- [2] J. Slay, M. Miller, Lessons learned from the Maroochy water breach, *Critical Infrastructure Protection* 253 (2007) 73–82.
- [3] S. Kuvshinkova, SQL Slammer worm lessons learned for consideration by the electricity sector, North American Electric Reliability Council.
- [4] J. P. Farwell, R. Rohozinski, Stuxnet and the future of cyber war, *Survival* 53 (1) (2011) 23–40.
- [5] Y. Liu, M. K. Reiter, P. Ning, False data injection attacks against state estimation in electric power grids, in: *ACM Conference on Computer and Communications Security*, Chicago, IL, USA, 2009, pp. 21–32.

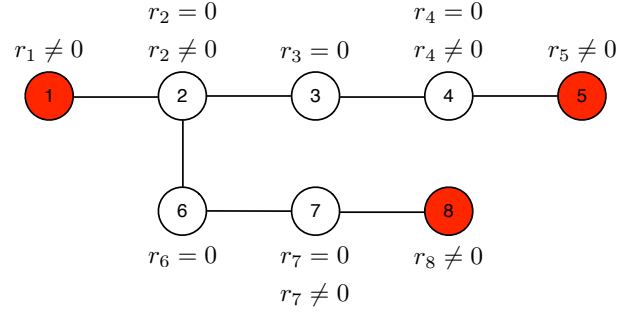


Figure B.4: An example of regional identification with 8 regions. The zero/nonzero pattern of the residuals computed by the control centers is reported: compromised regions {1, 5, 8} have nonzero residuals; non-compromised regions {3, 6} have zero residuals, because neighboring regions {2, 4, 7} are not compromised; regions 2, 4, and 7 may have zero or nonzero residuals, depending on the strategy of the attacker in region 1, 5, and 8, respectively. Non-compromised regions {2, 3, 4, 6, 7} are identified by Algorithm 1 as those regions with zero residual, and those with zero and nonzero neighboring residuals. Compromised regions {1, 5, 8} are identified by exclusion.

- [6] A. Teixeira, S. Amin, H. Sandberg, K. H. Johansson, S. Sastry, Cyber security analysis of state estimators in electric power systems, in: *IEEE Conf. on Decision and Control*, Atlanta, GA, USA, 2010, pp. 5991–5998.
- [7] R. Smith, A decoupled feedback structure for covertly appropriating network control systems, in: *IFAC World Congress*, Milan, Italy, 2011, pp. 90–95.
- [8] L. Xie, Y. Mo, B. Sinopoli, False data injection attacks in electricity markets, in: *IEEE Int. Conf. on Smart Grid Communications*, Gaithersburg, MD, USA, 2010, pp. 226–231.
- [9] F. Dörfler, F. Pasqualetti, F. Bullo, Continuous-time distributed observers with discrete communication, *IEEE Journal of Selected Topics in Signal Processing* 7 (2) (2013) 296–304.
- [10] G. Basile, G. Marro, *Controlled and Conditioned Invariants in Linear System Theory*, Prentice Hall, 1991.
- [11] F. Hamza, P. Tabuada, S. Diggavi, Secure state-estimation for dynamical systems under active adversaries, in: *Allerton Conf. on Communications, Control and Computing*, Monticello, IL, USA, 2011, pp. 337–344.
- [12] Y. Shoukry, P. Tabuada, Event-triggered state observers for sparse sensor noise/attacks, *arXiv preprint arXiv:1309.3511*.
- [13] T. Geerts, Invariant subspaces and invertibility properties for singular systems: The general case, *Linear Algebra and its Applications* 183 (1993) 61–88.
- [14] W. M. Wonham, *Linear Multivariable Control: A Geometric Approach*, 3rd Edition, Springer, 1985.
- [15] E. J. Candes, T. Tao, Decoding by linear programming, *IEEE Transactions on Information Theory* 51 (12) (2005) 4203–4215.
- [16] F. Pasqualetti, A. Bicchi, F. Bullo, Consensus computation in unreliable networks: A system theoretic approach, *IEEE Transactions on Automatic Control* 57 (1) (2012) 90–104.
- [17] M. Saif, Y. Guan, Decentralized state estimation in large-scale interconnected dynamical systems, *Automatica* 28 (1) (1992) 215–219.
- [18] F. J. Bejarano, T. Floquet, W. Perruquetti, G. Zheng, Observability and detectability analysis of singular linear systems with unknown inputs, in: *IEEE Conf. on Decision and Control and European Control Conference*, Orlando, FL, USA, 2011, pp. 4005–4010.
- [19] F. L. Lewis, Geometric design techniques for observers in singular systems, *Automatica* 26 (2) (1990) 411–415.