

Characterization of Model-Based Detectors for CPS Sensor Faults/Attacks

Carlos Murguía and Justin Ruths

Abstract—A vector-valued model-based cumulative sum (CUSUM) procedure is proposed for identifying faulty/falsified sensor measurements. First, given the system dynamics, we derive tools for tuning the CUSUM procedure in the fault/attack free case to fulfill a desired detection performance (in terms of false alarm rate). We use the widely-used chi-squared fault/attack detection procedure as a benchmark to compare the performance of the CUSUM. In particular, we characterize the state degradation that a class of attacks can induce to the system while enforcing that the detectors (CUSUM and chi-squared) do not raise alarms. In doing so, we find the upper bound of state degradation that is possible by an undetected attacker. We quantify the advantage of using a dynamic detector (CUSUM), which leverages the history of the state, over a static detector (chi-squared) which uses a single measurement at a time. Simulations of a chemical reactor with heat exchanger are presented to illustrate the performance of our tools.

Index Terms—Cyber Physical Systems, Model-based fault/attack detection, Security, CUSUM, Chi-squared.

I. INTRODUCTION

During the past half-century, scientific and technological advances have greatly improved the performance of control systems. From heating/cooling devices in our homes, to cruise-control in our cars, to robotics in manufacturing centers. However, these new technologies have also led to vulnerabilities of some of our most critical infrastructures—e.g., power, water, transportation. Advances in communication and computing power have given rise to adversaries with enhanced and adaptive capabilities. Depending on attacker resources and system defenses, attackers may deteriorate the functionality of systems even while remaining undetected. Therefore, designing efficient fault/attack detection schemes and attack-robust control systems is of key importance for guaranteeing the safety and proper operation of critical systems. Tools from sequential analysis and fault detection have to be adapted to deal with the systematic, strategic, and persistent nature of attacks. These new challenges have attracted the attention of many researchers in the control and computer science communities [1]-[10]. Lately, there has been increasing interest in studying *systems performance degradation* induced by attacks

that remain hidden or undetected by detection procedures [1]-[2],[5]-[6],[10]. Quantifying the system degradation provides a *measure of impact* to assess the performance of control structures, estimation schemes, and detection procedures against this class of intelligent attacks. For instance, in [5]-[6], for arbitrary detection procedures, the authors quantify how much the attacker can deviate the estimate of the state from its attack-free values while remaining stealthy. They characterize stealthiness of attacked sequences using the *Kullback-Leibler Divergence* [11] between the attack-free and the attacked sequence. In the same spirit, the authors in [1]-[2] study how attacks propagate through the control structure to degrade the system dynamics while remaining *undetected* by the detection mechanism. In particular, the authors in [1] characterize undetectability (for a class of deterministic LTI systems) as the ability of attackers to excite *only* the zero dynamics of the system [12] (making its effect undetectable from output measurements). In [2],[10],[13], the authors propose a notion of stealthiness by attacks that do not change the alarm rate of the detector by more than a small amount (thus making it hard for the operator to distinguish between an attack-free and an attacked system, i.e., these attacks remain *hidden* from the detector). As a measure of impact, they characterize the reachable sets that these hidden attacks can induce to the system.

Most of the current work on security of control systems has focused on *static* detection procedures (either bad-data or chi-squared detectors), which identify anomalies based on a single measurement at a time [1]-[4],[10]. There is only a small amount of literature considering the use of dynamic change detection procedures such as the Sequential Probability Ratio Test (SPRT) or the Cumulative Sum (CUSUM) [14], which employ measurement history, in the context of security of Cyber-Physical Systems (CPS) [7],[15]-[16]. Dynamic detectors present an appealing alternative to the aforementioned static procedures. Using measurement history provides extra degrees of freedom for improving the performance of our fault/attack detection strategies; in particular, against low amplitude persistent attacks [16].

This paper addresses, for Linear Time-Invariant (LTI) systems subject to sensor/actuator noise, the problem of characterizing CUSUM dynamic and chi-squared static detectors in terms of false alarm rates and performance degradation under a class of attacks. Standard Kalman filters are proposed to estimate the state of the physical process. Both detectors employ a distance measure that is a quadratic function of the residual (the error between sensor measurements and the estimated outputs). In the chi-squared procedure, at each time

Carlos Murguía is with the Center for Research in Cyber Security (iTrust), Engineering Systems and Design (ESD) Pillar, Singapore University of Technology and Design, Singapore. E-mail: murguia_rendon@sutd.edu.sg.

Justin Ruths is with the Departments of Mechanical and Systems Engineering, University of Texas at Dallas, USA, e-mail: jruths@utdallas.edu.

This work was supported by the National Research Foundation (NRF), Prime Minister's Office, Singapore, under its National Cybersecurity R&D Programme (Award No. NRF2014NCR-NCR001-40) and administered by the National Cybersecurity R&D Directorate.

Manuscript received May 5, 2017; revised May 5, 2017.

instant if the distance measure is larger than a threshold, an alarm is raised, indicating a possible compromised sensor. In the CUSUM procedure, the distance measure values are accumulated over time, and if this accumulated value is greater than expected an alarm is triggered. Fundamentally a detector aims to properly raise an alarm when a fault/attack happens and not raise an alarm when there is no fault/attack. Deviation from this ideal performance is captured by false positives (alarms are raised when there are no faults/attacks), also called false alarms, and false negatives (a fault/attack happens, but no alarm is raised). Although minimizing both false positives and false negatives is best, often they must be traded off based on which is more tolerable. In this context, detectors with high sensitivity would have high rates of false positives in favor of low rates of false negatives (and vice versa).

In order to provide an equitable comparison between detectors (e.g., static versus dynamic), we first require the ability to tune each type of detector to a similar level of sensitivity. To-date, however, there is not a complete characterization of how features of the system (e.g., system matrices, control/estimator gains, noise, sampling) affect the selection of the CUSUM parameters to achieve a desired sensitivity, quantified by the rate of false alarms. Our first contribution in this paper is to provide systematic tools to tune the CUSUM detector, and for completeness, also the chi-squared (static) detector, in the fault/attack free case based on the system dynamics, the Kalman filter, the stochastic properties of the distance measure, and a desired false alarm rate. In particular, sufficient conditions for mean square boundedness of the CUSUM sequence are derived when it is driven by a quadratic form of the residual. Then, using a Markov chain approximation of the CUSUM sequence, we give a procedure for selecting the *decision threshold* such that a desired false alarm rate is satisfied.

Second, for a class of *zero-alarm attacks* (attacks that prevent the detector from raising alarms), we characterize the impact of the attack sequence on the system dynamics when the vector-valued CUSUM and chi-squared detectors are deployed for attack detection. From an empirical point of view, zero-alarm attacks have been actively used to assess the resilience of dynamical systems against attacks [7],[16],[17]. zero-alarm attacks provide a simple, deterministic, yet representative class of attacks that can be easily scaled to systems with different dynamics and properties; thus making them a good choice for assessing the performance of attack detectors in terms of state degradation.

In our preliminary work [8], we have started analyzing these ideas. The contributions of this manuscript with respect to [8] are the following: A comprehensive and complete exposition of all the results and methodologies; the main results have been revised and improved and the corresponding proofs (which are not given in our preliminary work) are included in this paper; we formulate an all-new measure for attack degradation centered around the concept of input-to-state stability; and a benchmark simulation experiment used in the fault-detection literature [18],[19] (a chemical reactor with heat exchanger) is presented to illustrate the performance of our tools.

A. Notation

Throughout this paper, the following notation is used: the symbol \mathbb{R} stands for the real numbers, $\mathbb{R}_{>0}$ ($\mathbb{R}_{\geq 0}$) denotes the set of positive (non-negative) real numbers. The symbol \mathbb{N} stands for the set of natural numbers. The Euclidian norm in \mathbb{R}^n is denoted by $\|x\|$, $\|x\|^2 = x^T x$, where T denotes transposition. The induced norm of a matrix $A \in \mathbb{R}^{n \times n}$, denoted by $\|A\|$, is defined as $\|A\| = \max_{x \in \mathbb{R}^n, \|x\|=1} \|Ax\|$. The $n \times n$ identity matrix is denoted by I_n or simply I if no confusion can arise. Similarly, $n \times m$ matrices composed of only ones and only zeros are denoted by $\mathbf{1}_{n \times m}$ and $\mathbf{0}_{n \times m}$, respectively, or simply $\mathbf{1}$ and $\mathbf{0}$ when their dimensions are clear. If a quadratic form $x^T P x$ with a symmetric matrix $P = P^T$ is positive definite (semidefinite), then P is called positive definite (semidefinite). For positive definite (semidefinite) matrices, we use the notation $P > 0$ ($P \geq 0$); moreover, $P > Q$ ($P \geq Q$) means that the matrix $P - Q$ is positive definite (semidefinite). The spectrum of a matrix A is denoted by $\text{spec}[A]$, $\text{tr}[A]$ denotes its trace, and $\rho[A]$ is its spectral radius. The notation $\lambda_{\min}[A]$ ($\lambda_{\max}[A]$) stands for the smallest (largest) eigenvalue of the square matrix A . The notation $E[x]$ stands for the expected value of x and $E_y[x]$ denotes the expected value of x conditional to y . The variance of a random variable x is denoted by $\text{var}[x]$. The notation $\text{pr}[\cdot]$ denotes probability and $x \sim \mathcal{N}(\mu, \Sigma)$ means that $x \in \mathbb{R}^n$ is a vector-valued normally distributed random variable with mean $\mu \in \mathbb{R}^n$ and covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$. For simplicity of notation, we often suppress the explicit dependence of time t .

II. SYSTEM DESCRIPTION & ATTACK DETECTION

We study LTI stochastic systems of the form:

$$\begin{cases} x(t_{k+1}) = Fx(t_k) + Gu(t_k) + v(t_k), \\ y(t_k) = Cx(t_k) + \eta(t_k), \end{cases} \quad (1)$$

with sampling time-instants $t_k, k \in \mathbb{N}$, state $x \in \mathbb{R}^n$, measured output $y \in \mathbb{R}^m$, control input $u \in \mathbb{R}^l$, matrices F, G , and C of appropriate dimensions, and i.i.d. multivariate zero-mean Gaussian noises $v \in \mathbb{R}^n$ and $\eta \in \mathbb{R}^m$ with covariance matrices $R_1 \in \mathbb{R}^{n \times n}$, $R_1 \geq 0$ and $R_2 \in \mathbb{R}^{m \times m}$, $R_2 \geq 0$, respectively. The initial state $x(t_1)$ is assumed to be a Gaussian random vector with covariance matrix $R_0 \in \mathbb{R}^{n \times n}$, $R_0 \geq 0$. The processes $v(t_k), k \in \mathbb{N}$ and $\eta(t_k), k \in \mathbb{N}$ and the initial condition $x(t_1)$ are mutually independent. It is assumed that (F, G) is stabilizable and (F, C) is detectable. At the time-instants $t_k, k \in \mathbb{N}$, the output of the process $y(t_k)$ is sampled and transmitted over a communication network. The received output $\bar{y}(t_k)$ is used to compute control actions $u(t_k)$ which are sent back to the process, see Fig. 1. The complete control-loop is assumed to be performed instantaneously, i.e., the sampling, transmission, and arrival time-instants are equal. In this paper, we focus on attacks on sensor measurements. That is, in between transmission and reception of sensor data, an attacker may replace the signals coming from the sensors to the controller¹, see Fig. 1. After each transmission and reception,

¹Such an attack can also be accomplished by installing malware on the controller equipment, in which case true measurements reach the controller, but are manipulated before they are used.

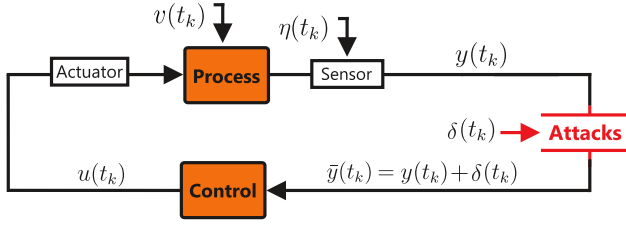


Fig. 1: Cyber-physical system under sensor attacks.

the attacked output \bar{y} takes the form:

$$\bar{y}(t_k) := y(t_k) + \delta(t_k) = Cx(t_k) + \eta(t_k) + \delta(t_k), \quad (2)$$

where $\delta(t_k) \in \mathbb{R}^m$ denotes *additive sensor attacks/faults*. Denote $x_k := x(t_k)$, $u_k := u(t_k)$, $v_k := v(t_k)$, $\bar{y}_k := \bar{y}(t_k)$, $\eta_k := \eta(t_k)$, and $\delta_k := \delta(t_k)$. Using this new notation, the attacked system is written in the following compact form:

$$\begin{cases} x_{k+1} = Fx_k + Gu_k + v_k, \\ \bar{y}_k = Cx_k + \eta_k + \delta_k. \end{cases} \quad (3)$$

Remark 1 If the stochastic processes $v_c(t_k)$ and $\eta(t_k)$ are non-Gaussian, using spectral factorization [20], [21], we could rewrite them as output signals coming from linear filters, say $G_1(q)$ and $G_2(q)$, with Gaussian stochastic processes as inputs, say $w_1(t_k)$ and $w_2(t_k)$; that is, $v_c(t_k) = G_1(q)w_1(t_k)$ and $\eta(t_k) = G_2(q)w_2(t_k)$, where q denotes the forward-shift operator. Then, by extending the system dynamics with the filters and considering the non-Gaussian noises $v_c(t)$ and $\eta(t_k)$ as new states, the extended system is written as a LTI system perturbed by Gaussian noise, see, for instance, [20], [21] for details.

A. Steady state Kalman filter (attack/fault free case)

To estimate the state of the process, a one step-ahead estimator with the following structure is proposed:

$$\hat{x}_{k+1} = F\hat{x}_k + Gu_k + L_k(\bar{y}_k - C\hat{x}_k), \quad (4)$$

with estimated state $\hat{x}_k \in \mathbb{R}^n$, $\hat{x}_1 = E[x(t_1)]$, and gain matrix $L_k \in \mathbb{R}^{n \times m}$. Define the estimation error $e_k := x_k - \hat{x}_k$. The matrix L_k is designed to minimize the covariance matrix $P_k := E[e_k e_k^T]$ in the absence of attacks. Given the discrete-time dynamics (3) and the estimator (4), the estimation error is governed by the difference equation:

$$e_{k+1} = (F - L_k C)e_k - L_k \eta_k - L_k \delta_k + v_k. \quad (5)$$

If the pair (F, C) is detectable, the covariance matrix converges to steady state in the sense that, in the attack-free case, $\lim_{k \rightarrow \infty} P_k = P$ exists [22]. Let $\delta_k = \mathbf{0}$; then, from (5), the mean value of e_k is given by

$$E[e_{k+1}] = (F - L_k C)E[e_k]. \quad (6)$$

Because $\hat{x}_1 = E[x(t_1)]$, the mean value of the estimation error equals $\mathbf{0}_{n \times 1}$ independent of L_k . We assume that the system has reached steady state before an attack occurs. Then, the estimation of the random sequence x_k , $k \in \mathbb{N}$ can be obtained by the estimator (4) with P_k and L_k in steady state. It can be verified that, if $CPC^T + R_2$ is positive definite (a standard

assumption that guarantees that the Kalman filter converges), the estimator gain:

$$L_k = L := (FPC^T)(R_2 + CPC^T)^{-1}, \quad (7)$$

leads to the minimal steady state covariance matrix P , with P given by the solution of the algebraic Riccati equation:

$$FPF^T - P + R_1 = FPC^T(R_2 + CPC^T)^{-1}CPF^T. \quad (8)$$

The reconstruction method given by (4)-(8) is referred to as the *steady state Kalman filter*, cf. [22].

Remark 2 It is well known that, if the noise sequences v_k and η_k are Gaussian, the Kalman filter (4)-(8) gives the best estimate \hat{x}_{k+1} of the state x_{k+1} (in terms of minimum-mean-square estimation error) from noisy measurements. Moreover, if the noise is not Gaussian, the Kalman filter is the best linear estimator; although there may exist nonlinear estimators with better performance, cf. [23].

B. Residuals and hypothesis testing

Attacks can be regarded as *induced faults* in the system. Then, it is reasonable to use existing fault detection techniques to identify sensor attacks. The main idea behind fault detection theory is the use of an *estimator* to forecast the evolution of the system in the absence of faults. This prediction is compared with the actual measurements from the sensors. If the difference between what it is measured and the estimation (often referred to as *residual*) is larger than expected, there might be a fault in the system. Although the notion of residuals and model-based detectors is now routine in the fault detection literature, the primary focus has been on detecting and isolating faults with *specific structures* (e.g., constant biases in sensor measurements or random faults in sensors and actuators following specific distributions). Now, in the context of an intelligent adversarial attacker, new challenges arise to understand the effect that an intruder can have on the system given the dynamics, the estimator, and the detector structure. In this work, we use the steady state Kalman filter introduced in the previous section as our estimator.

Consider the discrete-time process dynamics (3), the steady state Kalman filter (4)-(8), and the corresponding error difference equation (5). Define the residual sequence r_k , $k \in \mathbb{N}$ as

$$r_k := \bar{y}_k - C\hat{x}_k = Ce_k + \eta_k + \delta_k. \quad (9)$$

Then, r_k evolves according to the difference equation:

$$\begin{cases} e_{k+1} = (F - LC)e_k - L\eta_k + v_k - L\delta_k, \\ r_k = Ce_k + \eta_k + \delta_k. \end{cases} \quad (10)$$

If there are no faults/attacks, the mean of the residual is

$$E[r_{k+1}] = CE[e_{k+1}] + E[\eta_{k+1}] = \mathbf{0}_{m \times 1}, \quad (11)$$

and the covariance matrix is given by

$$E[r_{k+1}r_{k+1}^T] = CPC^T + R_2 =: \Sigma \in \mathbb{R}^{m \times m}. \quad (12)$$

For this residual, we identify two hypothesis to be tested: \mathcal{H}_0 the *normal mode* (no faults/attacks) and \mathcal{H}_1 the *faulty mode* (with faults/attacks). Under the normal mode, the statistics of the residual are:

$$\mathcal{H}_0 : \begin{cases} E[r_k] = \mathbf{0}_{m \times 1}, \\ E[r_k r_k^T] = \Sigma. \end{cases} \quad (13)$$

Therefore, when an fault/attack occurs in the system (\mathcal{H}_1), we expect that the statistics of the residual are different from the normal mode, i.e.,

$$\mathcal{H}_1 : \begin{cases} E[r_k] \neq \mathbf{0}_{m \times 1}, \text{ or} \\ E[r_k r_k^T] \neq \Sigma. \end{cases} \quad (14)$$

There exist many well-known hypothesis testing techniques which may be used to examine the residual and subsequently detect faults/attacks. For instance, Sequential Probability Ratio Testing (SPRT) [24], [25], Cumulative Sum (CUSUM) [14], [26], Generalized Likelihood Ratio (GLR) testing [27], Compound Scalar Testing (CST) [28], etc. Each of these techniques has its own advantages and disadvantages depending on the scenario. The most utilized and powerful one is, arguably the SPRT, which minimizes the time to reach a decision for given probabilities of false detection (i.e., declaring \mathcal{H}_1 when it is actually \mathcal{H}_0). In this manuscript, we mainly focus on the CUSUM procedure which is a version of SPRT that permits repeated detection [26]. However, for comparison, we also present results about a particular case of CST, namely the so-called *chi-squared* change detection procedure.

C. Distance measures and CUSUM procedure

Change detection theory was founded by Wald in 1947 when his book "*Sequential Analysis*" was published and the SPRT was first introduced. Subsequently, the CUSUM procedure [26] was proposed by Page to detect changes in the mean of random variables by testing a weighted sum of the last few observations, i.e., a moving average. As Page pointed out, the CUSUM is equivalent to a repeated SPRT in which the test is restarted once a change has been detected. The input to the CUSUM procedure is a *distance measure* $z_k \in \mathbb{R}$, $k \in \mathbb{N}$, i.e., a measure of how deviated the estimator is from the sensor measurements. We propose the quadratic distance measure

$$z_k := r_k^T \Sigma^{-1} r_k, \quad (15)$$

where r_k and Σ are the residual sequence and its covariance matrix defined in (9) and (12), respectively. If there are no attacks, $E[r_k] = \mathbf{0}$ and $E[r_k r_k^T] = \Sigma$; it follows that

$$\begin{cases} E[z_k] = \text{tr}[\Sigma^{-1} \Sigma] + E[r_k]^T \Sigma^{-1} E[r_k] \\ = m, \\ \text{var}[z_k] = 2\text{tr}[\Sigma^{-1} \Sigma \Sigma^{-1} \Sigma] + 4E[r_k]^T \Sigma^{-1} \Sigma \Sigma^{-1} E[r_k] \\ = 2m, \end{cases} \quad (16)$$

see, e.g., [11] for details. Moreover, since $r_k \sim \mathcal{N}(\mathbf{0}, \Sigma)$, then $z_k = r_k^T \Sigma^{-1} r_k$ follows a chi-squared distribution with m degrees of freedom, cf. [11]. Other options are based on *likelihood ratios*. In this case, instead of directly using the sequence z_k to drive the CUSUM procedure, the *log-likelihood ratio* $\Lambda_k(z_k)$ between the two hypotheses is employed:

$$\Lambda_k(z_k) := \log \frac{f_{z_k}^1(z|\mathcal{H}_1)}{f_{z_k}^0(z|\mathcal{H}_0)}, \quad (17)$$

where $f_{z_k}^j(z|\mathcal{H}_j)$ denotes the Probability Density Function (PDF) of the distance measure z_k , $k \in \mathbb{N}$ under \mathcal{H}_j , $j = \{0, 1\}$. A problem to address when using log-likelihood ratios for detecting attacks or unstructured faults is the fact that the PDF of the faulty sequence $f_1(z_k|\mathcal{H}_1)$ is unknown.

Actually, in the case of attacks, the adversary may induce any arbitrary (and possibly) non-stationary sequence z_k . Assuming the statistical properties of the attack sequences may limit our ability to detect a wide range of attacks [7].

The CUSUM procedure of Page driven by the distance measure z_k is defined as follows.

CUSUM:

$$\begin{cases} S_1 = 0, \\ S_k = \max(0, S_{k-1} + z_k - b), \text{ if } S_{k-1} \leq \tau, \\ S_k = 0 \text{ and } \tilde{k} = k - 1, \text{ if } S_{k-1} > \tau. \end{cases} \quad (18)$$

Design parameters: bias $b \in \mathbb{R}_{>0}$ and threshold $\tau \in \mathbb{R}_{>0}$.

Output: alarm time(s) \tilde{k} .

The idea is that the test sequence S_k accumulates the distance measure z_k and alarms are triggered when S_k exceeds the threshold τ . The test is reset to zero each time S_k becomes negative or larger than τ . If z_k is an independent non-negative sequence (which is our case) and b is not sufficiently large, the CUSUM sequence S_k grows unbounded until the threshold τ is reached, no matter how large τ is set. In order to prevent these drifts, inevitably leading to false alarms, the bias b must be selected properly based on the statistical properties of the distance measure. Once the the bias is chosen, the threshold τ must be selected to fulfill a required false alarm rate \mathcal{A}^* (see Section III-B).

III. CUSUM-TUNING

To enhance the performance of the CUSUM procedure, the bias b and the threshold τ must be selected appropriately. We have already mentioned that too small a bias can lead to inevitable growth of the CUSUM test sequence. At the same time, too large a bias may hide the effect of faults/attacks. In what follows, we provide tools for selecting these parameters given the statistical properties of the distance measure z_k introduced in (16). In particular, we provide sufficient conditions on the bias b such that, *in the absence of faults/attacks*, the sequence S_k of the CUSUM remains bounded (independent of the reset due to τ) in mean-squared sense. This is important to avoid false alarms due to the inherent divergence of S_k . Subsequently, we characterize the false alarm rate of the CUSUM in terms of b and τ given a *desired* false alarm rate.

A. Boundedness

First, we introduce the following concept of boundedness of stochastic processes, cf. [29],[30], followed by sufficient conditions for boundedness of the CUSUM sequence.

Definition 1 *The sequence S_k , $k \in \mathbb{N}$ is said to be bounded in mean square, if*

$$\sup_{k \in \mathbb{N}} E_{S_1} [S_k^2] < \infty,$$

is satisfied, i.e., the second moment of S_k is finite.

Theorem 1 *Consider the discrete-time process (3) and the steady state Kalman filter (4)-(8). Assume that there are no attacks to the system, i.e., $\delta_k = \mathbf{0}$. Let the CUSUM (18)*

with bias $b \in \mathbb{R}_{>0}$ and threshold $\tau \in \mathbb{R}_{>0}$ be driven by the distance measure $z_k = r_k^T \Sigma^{-1} r_k$, $k \in \mathbb{N}$ with residual sequence $r_k \sim \mathcal{N}(\mathbf{0}, \Sigma)$, $k \in \mathbb{N}$. Then, if the bias is set larger than the number of measurements, $b > \bar{b} := m$, the CUSUM sequence S_k , $k \in \mathbb{N}$ is bounded in mean square sense independent of the threshold τ .

The proof of Theorem 1 is presented in the appendix.

Remark 3 Notice that boundedness of the first moment follows from boundedness of the second moment and Jensen's inequality [11]. Then, $b > \bar{b} = m$ implies that the expected value $E_{S_1}[S_k]$, $k \in \mathbb{N}$ is finite.

The result stated in Theorem 1 implies that for $b > \bar{b}$, the second moment (and hence the first) of the sequence S_k , $k \in \mathbb{N}$ does not diverge. Consequently, we avoid false alarms due to intrinsic growth of the CUSUM sequence. Note that if the bias b is selected greater than but close to \bar{b} , small changes in the distance measure z_k would lead to divergence of S_k . Therefore, the smaller the bias, the higher the sensitivity against changes in (or uncertain characterization of) the residual signals.

B. False Alarms

Once the bias is selected such that boundedness of the second moment $E[S_k^2]$ is guaranteed, the next step is to select the threshold τ to fulfill a desired false alarm rate. The occurrence of an alarm in the CUSUM when there are no faults/attacks to the CPS is referred to as a false alarm. Operators need to tune this false alarm rate depending on the application. To do this, the threshold τ must be selected to fulfill a *desired false alarm rate* \mathcal{A}^* . Let $\mathcal{A} \in [0, 1]$ denote the *false alarm rate* of the procedure defined as the expected proportion of observations which are false alarms, i.e., for the CUSUM procedure, $\mathcal{A} := \text{pr}[S_k \geq \tau]$, see [31] and [32]. Define the *run length* \mathcal{K} of the CUSUM (18) as the number of iterations needed such that $S_{\mathcal{K}} > \tau$ (without attacks):

$$\mathcal{K} := \min\{k \geq 1 : S_k > \tau\}. \quad (19)$$

The expected value $E[\mathcal{K}]$ of \mathcal{K} is known in the literature as the *Average Run Length* (ARL). The ARL is inversely proportional to the false alarm rate \mathcal{A} [32], [31], i.e.,

$$\mathcal{A} = 1/\text{ARL}. \quad (20)$$

Then, for a given $b > \bar{b}$, the problem of selecting τ to satisfy a desired false alarm rate \mathcal{A}^* can be reformulated as the problem of selecting τ such that

$$\text{ARL} = 1/\mathcal{A}^*. \quad (21)$$

To determine a pair (b, τ) satisfying (21), an expression for the $\text{ARL} = E[\mathcal{K}]$ is required but, in general, its exact evaluation is analytically intractable [33]. The problem of approximating the ARL for CUSUM procedures has been addressed by many authors during the last decades. For instance, the authors in [33]-[35] propose Wiener process approximations of the ARL using analogies between the CUSUM and the SPRT

for normally distributed distance measures. Although these techniques lead to explicit formulas for evaluating the ARL, the obtained approximations are often too conservative, see [34]-[35]. Accurate numerical methods have been proposed by, for instance, [36]-[39]. These methods rely on two main techniques, namely Markov chain and integral equation approaches. Both methods give accurate predictions of the ARL (see [36] for a comparison); however, we find the Markov chain approach more constructive and easier to implement. In this work, we use the result of Evans and Brook [37]. With this result, we outline a procedure for selecting the threshold τ given the bias b and a required false alarm rate \mathcal{A}^* .

For given $b > \bar{b}$ and some $\tau \in \mathbb{R}_{>0}$, consider the sequence S_k generated by the CUSUM procedure (18) driven by the distance measure $z_k = r_k^T \Sigma^{-1} r_k$, $k \in \mathbb{N}$. Given the recursive nature of the CUSUM procedure and independence of v_k and η_k , $k \in \mathbb{N}$, the sequence S_k forms a Markov chain taking values on the non-negative real line [40]. By discretizing the probability distribution of the distance measure, it is possible to subdivide the CUSUM sequence S_k into a finite set of partitions. The idea is to approximate the continuous scheme by a Markov chain having $N + 1$ states labeled as $\{E_0, E_1, \dots, E_N\}$, where E_N is absorbing. Then, the probability that the chain remains in the same state at the next step should correspond to the case when S_k does not change in value by more than a small amount, say $\frac{1}{2}\Delta_S$, i.e., the next distance measure z_k does not differ from the bias b by more than $\frac{1}{2}\Delta_S$. The constant Δ_S determines the width of the grouping interval involved in the discretization of the probability distribution of z_k . The interval width $\frac{1}{2}\Delta_S$ must be selected such that the probability of jumping from E_j , $j \in \{0, \dots, N-1\}$ to the absorbing state E_N is approximately equal to the probability that the CUSUM sequence S_k jumps beyond the threshold τ from a position $S_{k-1} \in (0, \tau)$ which corresponds approximately to the state E_j . This requirement is satisfied by taking

$$\Delta_S := \frac{2\tau}{2N-1}, \quad (22)$$

see [37] for details. Then, the transition probabilities from a starting state E_j , $j = 0, \dots, N-1$, can be determined from the probability distribution of $z_k - b = r_k^T \Sigma^{-1} r_k - b$, as:

$$\begin{aligned} \text{pr}(E_j \rightarrow E_0) &= \text{pr}(z_k - b \leq -j\Delta_S + \frac{1}{2}\Delta_S), \\ \text{pr}(E_j \rightarrow E_N) &= \text{pr}((N-j)\Delta_S - \frac{1}{2}\Delta_S < z_k - b), \\ \text{pr}(E_j \rightarrow E_\nu) &= \text{pr}(z_k - b \leq (\nu-j)\Delta_S + \frac{1}{2}\Delta_S) \\ &\quad - \text{pr}(z_k - b < (\nu-j)\Delta_S - \frac{1}{2}\Delta_S). \end{aligned}$$

Note that $\text{pr}(E_0 \rightarrow E_N) = \text{pr}(z_k - b > \tau)$. For given b and τ , the states $\{E_0, \dots, E_N\}$ and the above transition probabilities forms a Markov chain whose transition matrix can be constructed from the probability distribution of $z_k - b$. Denote $T_\chi := \text{pr}(z_k - b \leq \chi\Delta_S + \frac{1}{2}\Delta_S)$ and $p_\chi := \text{pr}(\chi\Delta_S - \frac{1}{2}\Delta_S < z_k - b \leq \chi\Delta_S + \frac{1}{2}\Delta_S)$. Then, the Markov

transition matrix $\mathcal{P} \in \mathbb{R}^{(N+1) \times (N+1)}$ is given by:

$$\mathcal{P} := \begin{pmatrix} T_0 & p_1 & p_2 & \dots & p_{N-1} & 1 - T_{N-1} \\ T_{-1} & p_0 & p_1 & \dots & p_{N-2} & 1 - T_{N-2} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ T_{-j} & p_{1-j} & p_{2-j} & \dots & p_{N-1-j} & 1 - T_{N-1-j} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ T_{1-N} & p_{2-N} & p_{3-N} & \dots & p_0 & 1 - T_0 \\ 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix}. \quad (23)$$

Since the state E_N is absorbing, the last row consists of zeros except for the last entry. To compute the transition probabilities T_χ and p_χ of \mathcal{P} , we need the Cumulative Distribution Function (CDF) of the shifted distance measure $z_k - b = r_k^T \Sigma^{-1} r_k - b$. If there are no attacks, $r_k \sim \mathcal{N}(\mathbf{0}, \Sigma)$; therefore, $z_k - b$ follows a shifted chi-squared distribution with CDF:

$$F_{z_k - b}(x) := \begin{cases} \mathbf{P}\left(\frac{m}{2}, \frac{x+b}{2}\right), & \text{for } x \geq -b, \\ 0, & \text{for } x < -b, \end{cases} \quad (24)$$

where $\mathbf{P}(\cdot, \cdot)$ denotes the regularized lower incomplete gamma function [11]. Then, the entries of the transition matrix are given by

$$\begin{cases} p_\chi = F_{z_k - b}\left(\chi \Delta_S + \frac{1}{2} \Delta_S\right) \\ \quad - F_{z_k - b}\left(\chi \Delta_S - \frac{1}{2} \Delta_S\right), \\ T_\chi = F_{z_k - b}\left(\chi \Delta_S + \frac{1}{2} \Delta_S\right). \end{cases} \quad (25)$$

Define the transformation $\mathcal{T} := (I_N \ \mathbf{0}_{N \times 1}) \in \mathbb{R}^{N \times (N+1)}$ and the matrix:

$$\mathcal{R} := \mathcal{T} \mathcal{P} \mathcal{T}^T \in \mathbb{R}^{N \times N}. \quad (26)$$

The matrix \mathcal{R} is known as the *fundamental matrix* associated with the Markov transition matrix \mathcal{P} . Note that

$$\mathcal{P} = \begin{pmatrix} \mathcal{R} & * \\ \mathbf{0}_{1 \times N} & 1 \end{pmatrix}.$$

Then, all entries of \mathcal{R} are non-negative and its row sums are less than one. Therefore, by Gershgorin circle theorem, the eigenvalues of \mathcal{R} satisfy: $1 > |\lambda_N| \geq \dots \geq |\lambda_1|$. It follows that $\rho[\mathcal{R}] < 1$, where $\rho[\cdot]$ denotes spectral radius; therefore, the matrix $(I_N - \mathcal{R})$ is invertible [41]. Next, having introduced the transition matrix \mathcal{P} of the approximated Markov chain and the fundamental matrix \mathcal{R} , we can compute an approximation $\tilde{\mathcal{A}}$ of the false alarm rate \mathcal{A} based on the result in [37], equation (20), and (22)-(26).

Theorem 2 *Assume that there are no attacks on the system and let the CUSUM (18) with bias $b > \bar{b} = m$ and threshold $\tau \in \mathbb{R}_{>0}$ be driven by the distance measure $z_k = r_k^T \Sigma^{-1} r_k$ with residual sequence $r_k \sim \mathcal{N}(\mathbf{0}, \Sigma)$, $k \in \mathbb{N}$. For a finite number of partitions $N \in \mathbb{N}$, consider the fundamental matrix \mathcal{R} , defined in (26), obtained from the transition matrix \mathcal{P} (22)-(25), and define*

$$\mu := (I_N - \mathcal{R})^{-1} \mathbf{1}_{N \times 1} = [\mu_1, \dots, \mu_N]^T. \quad (27)$$

Then, the false alarm rate $\mathcal{A} = 1/\text{ARL}$, is approximately given by $\tilde{\mathcal{A}} := \mu_1^{-1}$. Moreover, as $N \rightarrow \infty$, $\tilde{\mathcal{A}} \rightarrow \mathcal{A}$, i.e., $\lim_{N \rightarrow \infty} \tilde{\mathcal{A}} = \mathcal{A}$.

Proof: Consider the Markov chain \mathcal{M} given by the states $\{E_0, E_1, \dots, E_N\}$ and the transition matrix \mathcal{P} (22)-(25). Let $\tilde{\mathcal{K}} \in \mathbb{N}$ denote the number of iterations needed to reach the absorbing state E_N from E_0 . The random variable $\tilde{\mathcal{K}}$ follows a *discrete phase-type distribution* with $E[\tilde{\mathcal{K}}] =$

μ_1 and μ_1 as defined in (27), see [42]. By construction, \mathcal{M} is a finite state approximation of the continuous Markov chain formed by the CUSUM sequence $S_k \in \mathbb{R}_{\geq 0}$, $k \in \mathbb{N}$ driven by $z_k = r_k^T \Sigma^{-1} r_k$, $k \in \mathbb{N}$. It follows that $E[\tilde{\mathcal{K}}] = \mu_1 \approx E[\mathcal{K}] = \text{ARL}$ where \mathcal{K} denotes the *run length* of the CUSUM defined in (19). Then, from (20), we have that $\mathcal{A} = \text{ARL}^{-1} \approx \mu_1^{-1}$. Next, increasing the number of partitions N would reduce the width of the grouping interval Δ_S (22), such that, as $N \rightarrow \infty$, the Markov chain \mathcal{M} retrieves the continuous scheme given by the CUSUM sequence $S_k \in \mathbb{R}_{\geq 0}$, $k \in \mathbb{N}$, (18); therefore, $\mathcal{A} = 1/\lim_{N \rightarrow \infty} E[\tilde{\mathcal{K}}]$. ■

Remark 4 *Theorem 2 provides a tool for approximating the false alarm rate \mathcal{A} of the CUSUM procedure for given bias b and threshold τ . In particular, for a given $b > \bar{b}$, it provides a map $\mathcal{S} : \mathbb{R}_{>0} \rightarrow (0, 1)$ from the threshold τ to the approximated false alarm rate $\tilde{\mathcal{A}}$, i.e., $\tau \mapsto \mathcal{S}(\tau)$, $\tilde{\mathcal{A}} = \mathcal{S}(\tau)$. Given that $F_{z_k - b}(z)$ is a continuous function for all $z \in \mathbb{R}$, it can be proved that $\mathcal{S}(\tau)$ is continuous for all $\tau \in \mathbb{R}_{>0}$; then, simple bisection methods can be used to determine the threshold $\tau = \tau^* \in \mathbb{R}_{>0}$ required to satisfy $\tilde{\mathcal{A}} = \mathcal{S}(\tau^*) = \mathcal{A}^*$ for given $b > \bar{b}$.*

IV. CHI-SQUARED TUNING

The CUSUM approach to fault/attack detection offers an compelling alternative to the more popular chi-squared detector. Here, we use the chi-squared approach as a benchmark to compare the performance of the CUSUM. Consider again the residual sequence r_k , (10), and its covariance matrix Σ , (12). The chi-squared procedure is defined as follows:

Chi-squared procedure:

$$\text{If } z_k = r_k^T \Sigma^{-1} r_k > \alpha, \quad \tilde{k} = k. \quad (28)$$

Design parameter: threshold $\alpha \in \mathbb{R}_{>0}$.

Output: alarm time(s) \tilde{k} .

The idea is that alarms are triggered if z_k exceeds the threshold α . Similar to the CUSUM procedure, the parameter α is selected to satisfy a required false alarm rate \mathcal{A}^* .

Theorem 3 *Assume that there are no attacks on the system and consider the chi-squared procedure (28) with threshold $\alpha \in \mathbb{R}_{>0}$, $r_k \sim \mathcal{N}(\mathbf{0}, \Sigma)$. Let $\alpha = \alpha^* := 2\mathbf{P}^{-1}\left(\frac{m}{2}, 1 - \mathcal{A}^*\right)$, where $\mathbf{P}^{-1}(\cdot, \cdot)$ denotes the inverse regularized lower incomplete gamma function, then $\mathcal{A} = \mathcal{A}^*$.*

Proof: Let $\tilde{\mathcal{K}}$ denote the run length of the chi-squared procedure (28) defined as the number of iterations needed such that $k = \tilde{\mathcal{K}}$ implies $z_k > \alpha$ when there are no attacks. As with the CUSUM procedure, the Average Run Length is given by $\text{ARL} = E[\tilde{\mathcal{K}}]$ and the false alarm rate satisfies $\mathcal{A} = 1/\text{ARL}$. The random variable $\tilde{\mathcal{K}}$ follows a geometric distribution [11]; therefore, $\text{ARL} = E[\tilde{\mathcal{K}}] = 1/\text{pr}(z_k > \alpha)$ and $\mathcal{A} = \text{pr}(z_k > \alpha)$. Each element of the sequence z_k , $k \in \mathbb{N}$ is an i.i.d. random variable with CDF given by $\tilde{F}_{z_k}(x) = \mathbf{P}\left(\frac{m}{2}, \frac{x}{2}\right)$. Then, $\mathcal{A} = \text{pr}(z_k > \alpha) = 1 - \mathbf{P}\left(\frac{m}{2}, \frac{\alpha}{2}\right)$ and the result follows. ■

V. DETECTOR PERFORMANCE UNDER ZERO-ALARM ATTACKS

In this section, we assess the performance of the CUSUM procedure by quantifying the effect of the attack sequence δ_k on the estimation error when the CUSUM procedure is used to identify anomalies. To maintain an equitable comparison between detectors in this section, some assumption must be made about their false positive rate (false alarm rate) and false negative rate (the rate at which true attacks are not detected). Using the tools introduced in prior sections, we can calibrate the CUSUM and chi-squared detectors to have the same false alarm rate. Here, we consider a class of *zero-alarm attacks*, i.e., attack sequences that keep the detector from raising alarms. This implies that the entire attacked distribution (of z_k or S_k) is at or below the decision threshold, effectively maximizing the false negative rate (since the true positive rate is zero, i.e., the true attack is never detected). Zero-alarm attacks provide a concise quantification of attacker impact on the system performance. In particular, we characterize the estimation error deviation due to zero-alarm attacks. This serves as a useful proxy for the capabilities of the attacker due to the detection mechanism. To do this, using the notion of *input to state stability* [43]-[44], we derive upper bounds on the trajectories of the estimation error given the system dynamics, the attack sequence, and the CUSUM parameters. Furthermore, we compare the performance of the CUSUM against the chi-squared detector. In this paper, we have now characterized the rate of false alarms based on the type of detector and the system and detector parameters. We now close the loop on this analysis by identifying the impact that attackers can have while taking advantage of the detector structure.

A. Zero-alarm Attacks

Here, we quantify the damage that attacks may induce to the estimation error dynamics while enforcing that alarms are not raised by the detector. We assume that the attacker has perfect knowledge of the system dynamics, the Kalman filter, control inputs, measurements, and detection procedure (either CUSUM or chi-squared). It is further assumed that all the sensors can be compromised by the attacker at each time step (a worst-case scenario).

First, consider the chi-squared procedure (28) and write z_k in terms of the estimation error e_k :

$$z_k = (Ce_k + \eta_k + \delta_k)^T \Sigma^{-1} (Ce_k + \eta_k + \delta_k). \quad (29)$$

Because e_k and η_k have infinite support, to prevent z_k from going beyond the threshold α , the attack sequence δ_k must compensate for the term $Ce_k + \eta_k$. By assumption, the attacker has access to $y_k = Cx_k + \eta_k$ (real-time sensor measurements). Moreover, given its perfect knowledge of the Kalman filter, the adversary can compute the estimated output $C\hat{x}_k$ and then construct $y_k - C\hat{x}_k = Ce_k + \eta_k$. For a given chi-squared threshold α , define the sequence $\bar{\delta}_k^\alpha := \{\bar{\delta}_k^\alpha \in \mathbb{R}^m \mid (\bar{\delta}_k^\alpha)^T \bar{\delta}_k^\alpha \leq \alpha\}$; for instance, $\bar{\delta}_k^\alpha = [\sqrt{\frac{\alpha}{m}}, \sqrt{\frac{\alpha}{m}}, \dots, \sqrt{\frac{\alpha}{m}}]^T$ and $\bar{\delta}_k^\alpha = [\sqrt{\alpha}, 0, \dots, 0]^T$. Let $k = k^*$ denote the starting attack instant for some $k^* \geq 1$. Then, for $k \geq k^*$, it follows that

$$\delta_k = -Ce_k - \eta_k + \Sigma^{\frac{1}{2}} \bar{\delta}_k^\alpha \rightarrow z_k \leq \alpha, \quad (30)$$

where $\Sigma^{\frac{1}{2}}$ denotes the symmetric square root matrix of Σ , is a feasible attack sequence given the capabilities of the attacker. Sequences δ_k of the form (30) define a class of attacks that can be launched by the opponent while preventing the chi-squared detector from raising alarms, i.e., zero-alarm attacks. The estimation error dynamics under the attack (30) is given by

$$e_{k+1} = Fe_k - L\Sigma^{\frac{1}{2}} \bar{\delta}_k^\alpha + v_k, \quad k \geq k^*. \quad (31)$$

Remark 5 Note that if $\rho[F] > 1$, then $\|E[e_k]\|$ diverges to infinity as k grows for any nonstabilizing $\bar{\delta}_k^\alpha$ [22]. That is, zero-alarm attacks of the form (30) may destabilize the system if $\rho[F] > 1$. If $\rho[F] \leq 1$, then $\|E[e_k]\|$ may or may not diverge to infinity depending on algebraic and geometric multiplicities of the eigenvalues with unit modulus of F (a known fact from stability of LTI systems [22]).

Using the superposition principle of linear systems, the estimation error e_k can be written as $e_k = e_k^v + e_k^\delta$, where e_k^v denotes the part of e_k driven by noise and e_k^δ is the part driven by attacks. Using this new notation, we can write the dynamics (31) as follows:

$$e_{k+1}^v = Fe_k^v + v_k, \quad (32)$$

$$e_{k+1}^\delta = Fe_k^\delta - L\Sigma^{\frac{1}{2}} \bar{\delta}_k^\alpha, \quad k \geq k^*, \quad (33)$$

with $e_{k^*}^v = e_{k^*}$ and $e_{k^*}^\delta = \mathbf{0}$. Therefore, the contribution of zero-alarm attacks to e_k is solely determined by e_k^δ generated by (33). For a sequence $s_k \in \mathbb{R}^n$, $k \in \mathbb{N}$, let $s_{[k^*, k]}$ denote the truncation of s_k from k^* to k , i.e., $s_{[k^*, k]} := \{s_{k^*}, \dots, s_k\}$ and $\|s_{[k^*, k]}\| := \sup_{k^* \leq N \leq k} \|s_N\|$. For any matrix $A^{n \times n}$ such that $\rho[A] < 1$, let $\|\cdot\|_*$ denote some matrix norm satisfying $\|A\|_* < 1$ (such a norm always exists if $\rho[A] < 1$ [41]).

Proposition 1 Consider the process (3), the Kalman filter (4)-(8), and the chi-squared procedure (28) with threshold $\alpha \in \mathbb{R}_{>0}$. Assume $\rho[F] < 1$ and let $c \in \mathbb{R}_{>0}$ be some constant satisfying $\|F^k\| \leq c \|F^k\|_*$ for all $k \in \mathbb{N}$. Let the sensors be attacked by the sequence (30); then, for all $\bar{\delta}_{[k^*, k]}^\alpha$, $k > k^* \in \mathbb{N}$, the trajectories of (33) satisfy the inequalities:

$$\begin{cases} \|e_k^\delta\| \leq \gamma_k^{\chi^2} := \sqrt{\alpha} c \|L\Sigma^{\frac{1}{2}}\| \frac{1 - \|F\|_*^{k-k^*}}{1 - \|F\|_*}, \\ \lim_{k \rightarrow \infty} \|e_k^\delta\| \leq \bar{\gamma}^{\chi^2} := \frac{\sqrt{\alpha} c \|L\Sigma^{\frac{1}{2}}\|}{1 - \|F\|_*}. \end{cases} \quad (34)$$

Proof: The solution of (33), for $k > k^*$, is given by

$$e_k^\delta = - \sum_{i=0}^{k-1-k^*} F^i L\Sigma^{\frac{1}{2}} \bar{\delta}_{k-i-1}^\alpha,$$

it follows that

$$\|e_k^\delta\| \leq \|L\Sigma^{\frac{1}{2}}\| \sum_{i=0}^{k-1-k^*} \|F^i\| \|\bar{\delta}_{[k^*, k-1]}^\alpha\|, \quad k > k^*.$$

Because $\rho[F] < 1$, there exists a matrix norm, say $\|\cdot\|_*$, such that $\|F\|_* < 1$ (see Lemma 5.6.10 in [41]). Moreover, because all norms are equivalent in finite dimensional vector spaces [41], there exists a constant $c \in \mathbb{R}_{>0}$ satisfying $\|D\| \leq c \|D\|_*$ for all $D \in \mathbb{R}^{n \times n}$ [45]. It follows that

$$\begin{aligned} \|e_k^\delta\| &\leq c \|L\Sigma^{\frac{1}{2}}\| \sum_{i=0}^{k-1-k^*} \|F\|_*^i \|\bar{\delta}_{[k^*, k-1]}^\alpha\| \\ &\leq c \|L\Sigma^{\frac{1}{2}}\| \frac{1 - \|F\|_*^{k-k^*}}{1 - \|F\|_*} \|\bar{\delta}_{[k^*, k-1]}^\alpha\|, \end{aligned}$$

because $\sum_{i=0}^n \|F\|_*^i$ is a geometric series. By construction, for all $k \geq k^*$, $(\bar{\delta}_k^\alpha)^T \bar{\delta}_k^\alpha \leq \alpha$ which implies $\|\bar{\delta}_{[k^*, k-1]}^\alpha\| = \sup_{k^* \leq N \leq k-1} \|\bar{\delta}_N^\alpha\| \leq \sqrt{\alpha}$; therefore, the estimation error driven by attacks e_k^δ satisfies the inequalities in (34). ■

The effect of the attack sequence (30) on the upper bound of the estimation error (34) is determined by the sequence $\gamma_k^{\chi^2}$. The sequence $\gamma_k^{\chi^2}$ quantifies the impact of zero-alarm attacks when the chi-squared detector is used to detect anomalies, i.e., $\gamma_k^{\chi^2}$ gives a measure of the detector performance for the class of attacks in (30) in terms of estimation error deviation. The sequence $\gamma_k^{\chi^2}$ depends on the norm $\|\cdot\|_*$ which could be any matrix norm satisfying $\|F\|_* < 1$. If usual matrix norms (e.g., $\|\cdot\|_1$, $\|\cdot\|_2$, $\|\cdot\|_\infty$, etc.) do not satisfy this condition, in the proof of Lemma 5.6.10 in [41], the authors give a procedure for constructing such a norm provided that $\rho(F) < 1$. For given norm $\|\cdot\|_*$ satisfying $\|F\|_* < 1$, the constant c can be taken as $c = \inf\{c \in \mathbb{R}_{>0} : \|F^k\| - c \|F^k\|_* \leq 0, \forall k \in \mathbb{N}\}$ which can be obtained numerically.

Next, consider the CUSUM procedure and write (18) in terms of the estimation error e_k :

$$S_k = \max(0, S_{k-1} + \|\Sigma^{-\frac{1}{2}}(Ce_k + \eta_k + \delta_k)\|^2 - b), \quad (35)$$

if $S_{k-1} \leq \tau$; and $S_k = 0$, if $S_{k-1} > \tau$. As with the chi-squared procedure, we look for attack sequences that maintain the CUSUM statistic below the threshold τ preventing alarms to be raised. Let the attack start at some $k = k^* \geq 2$ and $S_{k^*-1} \leq \tau$, i.e., the attack does not start immediately after a false alarm. Define $\bar{\tau}_k := \{\bar{\tau}_k \in \mathbb{R}^m | \bar{\tau}_k^T \bar{\tau}_k \leq \tau + b - S_{k-1}\}$ and $\bar{\delta}_k^b := \{\bar{\delta}_k^b \in \mathbb{R}^m | (\bar{\delta}_k^b)^T \bar{\delta}_k^b \leq b\}$ for given threshold τ and bias b . Consider the attack sequence:

$$\delta_k = \begin{cases} -Ce_k - \eta_k + \Sigma^{\frac{1}{2}} \bar{\tau}_k, & k = k^*, \\ -Ce_k - \eta_k + \Sigma^{\frac{1}{2}} \bar{\delta}_k^b, & k > k^*. \end{cases} \quad (36)$$

It follows that $S_{k^*} = \max(0, S_{k^*-1} + \bar{\tau}_{k^*}^T \bar{\tau}_{k^*} - b) \leq \tau$, $S_{k^*+1} = \max(0, S_{k^*} + (\bar{\delta}_{k^*+1}^b)^T \bar{\delta}_{k^*+1}^b - b) \leq \tau$, and $S_{k^*+N} = \max(0, S_{k^*+N-1} + (\bar{\delta}_{k^*+1}^b)^T \bar{\delta}_{k^*+1}^b - b) \leq \tau$ for all $N \in \mathbb{N}$. That is, the class of attack sequences in (36) prevents the CUSUM procedure from raising alarms. Note that the attacker can only induce this sequence by exactly knowing S_{k^*-1} , i.e., the value of the CUSUM sequence one step before the attack. This is a strong assumption since it represents a real-time quantity that is not communicated over the communication network. Even if the opponent has access to the parameters of the CUSUM (b, τ) given the stochastic nature of the residual, the attacker would need to know the complete history of observations (from when the CUSUM was started) to be able to reconstruct S_{k^*-1} from data. This is an inherent security advantage in favor of the CUSUM over static detectors like the bad-data or chi-squared. Nevertheless, for evaluating the worst case scenario, we assume that the attacker has access to S_{k^*-1} . By construction, the estimation error dynamics under the attack sequence (36) is written as: $e_{k^*+1} = Fe_{k^*}^* - L\Sigma^{\frac{1}{2}}\bar{\tau}_{k^*} + v_{k^*}$, and, for $k > k^*$,

$$e_{k+1} = Fe_k - L\Sigma^{\frac{1}{2}}\bar{\delta}_k^b + v_k. \quad (37)$$

Note that (31) and (37) have the same dynamics but different initial condition. Therefore, we may expect upper bounds on $\|e_k\|$ similar to (34) obtained for the chi-squared. Again, we write e_k as $e_k = e_k^v + e_k^\delta$, where e_k^v denotes the part of e_k

driven by noise and e_k^δ is the part driven by attacks of (37). Using this new notation, we can write the dynamics (37) as:

$$e_{k+1}^v = Fe_k^v + v_k, \quad (38)$$

$$e_{k+1}^\delta = Fe_k^\delta - L\Sigma^{\frac{1}{2}}\bar{\delta}_k^b, \quad k > k^*, \quad (39)$$

with $e_{k^*}^\delta = \mathbf{0}$, $e_{k^*+1}^\delta = -L\Sigma^{\frac{1}{2}}\bar{\tau}_{k^*}$, $e_{k^*}^v = e_{k^*}$, and $e_{k^*+1}^v = Fe_{k^*} + v_{k^*}$.

Proposition 2 Consider the process (3), the Kalman filter (4)-(8), and the CUSUM procedure (18) with threshold $\tau \in \mathbb{R}_{>0}$ and bias $b > \bar{b} = m \in \mathbb{N}_{>0}$. Assume $\rho[F] < 1$ and let $c \in \mathbb{R}_{>0}$ be some constant satisfy $\|F^k\| \leq c \|F^k\|_*$ for all $k \in \mathbb{N}$. Let the sensors be attacked by the sequence (36); then, for all $\bar{\tau}_{k^*}$ and $\bar{\delta}_{[k^*, k]}^b$, $k > k^* \in \mathbb{N}$, the trajectories of (39) satisfy the inequalities:

$$\begin{cases} \|e_k^\delta\| \leq \gamma_k^{\text{CS}} := \sqrt{b}c \|L\Sigma^{\frac{1}{2}}\| \frac{1 - \|F\|_*^{k-k^*}}{1 - \|F\|_*} \\ \quad + c \|L\Sigma^{\frac{1}{2}}\bar{\tau}_{k^*}\| \|F\|_*^{k-k^*-1}, \\ \lim_{k \rightarrow \infty} \|e_k^\delta\| \leq \bar{\gamma}^{\text{CS}} := \frac{\sqrt{b}c \|L\Sigma^{\frac{1}{2}}\|}{1 - \|F\|_*}. \end{cases} \quad (40)$$

Proof: The solution of (39), for $k > k^* + 1$, is given by

$$e_k^\delta = -F^{k-k^*-1}L\Sigma^{\frac{1}{2}}\bar{\tau}_{k^*} - \sum_{i=0}^{k-2-k^*} F^i L\Sigma^{\frac{1}{2}}\bar{\delta}_{k-i-1}^b.$$

Then, following the same lines as in the proof of Proposition 1, we can write the following

$$\begin{aligned} \|e_k^\delta\| &\leq c \|L\Sigma^{\frac{1}{2}}\bar{\tau}_{k^*}\| \|F\|_*^{k-k^*-1} \\ &\quad + c \|L\Sigma^{\frac{1}{2}}\| \frac{1 - \|F\|_*^{k-k^*}}{1 - \|F\|_*} \|\bar{\delta}_{[k^*+1, k-1]}^b\|, \end{aligned}$$

By construction, for all $k > k^*$, $(\bar{\delta}_k^b)^T \bar{\delta}_k^b \leq b$ which implies $\|\bar{\delta}_{[k^*+1, k-1]}^b\| = \sup_{k^*+1 \leq N \leq k-1} \|\bar{\delta}_N^b\| \leq \sqrt{b}$; therefore, the trajectories of the estimation error dynamics (39) satisfy the inequalities in (40). ■

When considering the CUSUM, the effect of attacks of the form (36) on the upper bound of the estimation error is determined by the sequence γ_k^{CS} in (40). Note that, in steady state, the $(\bar{\tau}_{k^*})$ -dependent term in γ_k^{CS} decreases exponentially to zero. Then, the sequences $\gamma_k^{\chi^2}$ and γ_k^{CS} take constant values asymptotically. It follows that we can directly quantify the performance ratio (in terms of steady state deviation of $\|e_k^\delta\|$) between the two detectors. This is stated in the following corollary of Proposition 1 and Proposition 2.

Corollary 1 In steady state:

$$\lim_{k \rightarrow \infty} \frac{\gamma_k^{\chi^2}}{\gamma_k^{\text{CS}}} = \frac{\bar{\gamma}^{\chi^2}}{\bar{\gamma}^{\text{CS}}} = \sqrt{\frac{\alpha}{b}}. \quad (41)$$

B. Detector Comparison

For the case of zero-alarm attacks, we have derived upper bounds on e_k^δ for both the chi-squared and CUSUM procedures provided that $\rho[F] < 1$ (otherwise $\|e_k^\delta\|$ diverges under the class of zero-alarm attacks considered). To compare these results, we use the ratio of sequences $(\gamma_k^{\chi^2}/\gamma_k^{\text{CS}})$. From Corollary 1, $\lim_{k \rightarrow \infty} (\gamma_k^{\chi^2}/\gamma_k^{\text{CS}}) = \sqrt{\alpha/b}$. That is, in steady state, the $(\bar{\tau}_{k^*})$ -dependent term in γ_k^{CS} exponentially decreases

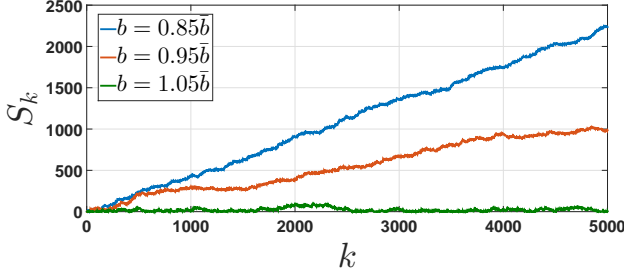


Fig. 2: CUSUM evolution for different values of bias b .

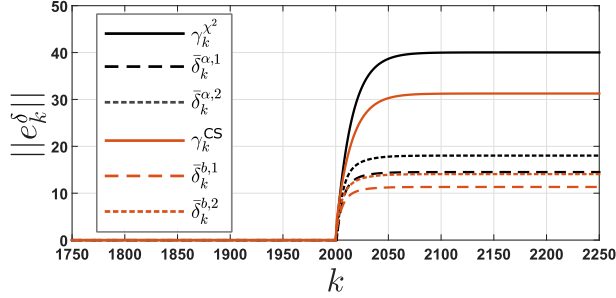


Fig. 3: Upper bounds $\gamma_k^{\chi^2}$ and γ_k^{CS} , and deviation of $\|e_k^\delta\|$ due to the zero-alarm attacks $(\bar{\delta}_k^{\alpha,1}, \bar{\delta}_k^{b,1})$ and $(\bar{\delta}_k^{\alpha,2}, \bar{\delta}_k^{b,2})$. Attacks are induced at $k = 2 \times 10^3$.

to zero and, if $b < \alpha$, under the same class of zero-alarm attacks, the CUSUM procedure leads to smaller steady state deviations on $\|e_k\|$ than the chi-squared procedure. In general, to increase the chances of attack detection, it is desired to select b as close as possible to \bar{b} in Theorem 1. It follows that $b \approx \bar{b} = m$. On the other hand, according to Theorem 3, α must be selected as $\alpha = \alpha^* = 2P^{-1}(\frac{m}{2}, 1 - \mathcal{A}^*)$ to fulfill a desired false alarm rate \mathcal{A}^* . In this case, we want to select \mathcal{A}^* close to zero, such that there are only a few false alarms. Let $\mathcal{A}^* \in \{0.01, 0.1\}$ and $m = 2$, i.e., false alarms between 1% and 10% and two dimensional outputs; then, $\alpha = 2P^{-1}(\frac{m}{2}, 1 - \mathcal{A}^*) \in [4.60, 9.21]$ and $b \approx \bar{b} = 2$. This implies that, in steady state, for the same class of attacks and $\mathcal{A}^* \in [0.01, 0.1]$, the chi-squared procedure leads to at least two times larger upper bounds than the CUSUM. Actually, for having $\alpha = b$ (which implies $\lim_{k \rightarrow \infty} (\gamma_k^{\chi^2} / \gamma_k^{\text{CS}}) = 1$), it is necessary to allow for a rate of $\mathcal{A}^* = 0.63$, which is high for practical purposes. For the CUSUM procedure, the threshold τ is selected to fulfill the desired \mathcal{A}^* . Given that there are no exact closed-form expressions to relate τ and \mathcal{A}^* (we provided a numeric approximation), it is not possible to exactly tell how large the τ would need be to satisfy \mathcal{A}^* . However, as mentioned, the contribution of $\bar{\tau}_{k^*}$ to e_k vanishes exponentially, i.e., independent of how large τ is, its contribution to γ_k^{CS} is zero in steady state.

VI. SIMULATION EXPERIMENTS

The authors in [18],[19] study the fault detection problem for a well stirred chemical reactor with heat exchanger. We

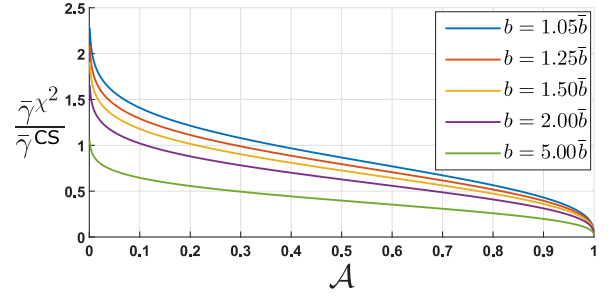


Fig. 4: Asymptotic ratio $(\bar{\gamma}^{\chi^2} / \bar{\gamma}^{\text{CS}})$ versus the false alarm rate \mathcal{A} for different values of CUSUM bias b .

use this system to demonstrate our results. The state, input, and output vectors of the considered reactor are:

$$\begin{cases} x(t) := \begin{pmatrix} C_0 \\ T_0 \\ T_w \\ T_m \end{pmatrix}, u(t) := \begin{pmatrix} C_u \\ T_u \\ T_{w,u} \end{pmatrix}, y(t) := \begin{pmatrix} C_0 \\ T_0 \\ T_w \end{pmatrix}, \end{cases}$$

where

$$\begin{cases} C_0 & : \text{Concentration of the chemical product,} \\ T_0 & : \text{Temperature of the product,} \\ T_w & : \text{Temperature of the jacket water of heat exchanger,} \\ T_m & : \text{Coolant temperature,} \\ C_u & : \text{Inlet concentration of reactant,} \\ T_u & : \text{Inlet temperature,} \\ T_{w,u} & : \text{Coolant water inlet temperature.} \end{cases}$$

We linearize the nonlinear model introduced in [18] about the origin $x(t) = \mathbf{0}_{4 \times 1}$ and then discretize it with sampling time $h = 0.05$. The resulting discrete-time linear system is given by (3)-(8) with matrices as given in (42). The original model in [18] does not consider sensor/actuator noise, we have included noise to increase the complexity of our simulation experiments. First, assume no attacks, i.e., $\delta_k = \mathbf{0}$, and consider the CUSUM procedure (18) with distance measure $z_k = r_k^T \Sigma^{-1} r_k$ and residual sequence (10). According to Theorem 1, the bias b must be selected larger than $\bar{b} = m = 3$ to ensure mean square boundedness of S_k independent of the threshold τ . Figure 2 depicts the evolution of the CUSUM for $b \in \{0.85\bar{b}, 0.95\bar{b}, 1.05\bar{b}_1\}$ and $k \in [1, 5000]$. For the purpose of illustrating this unbounded growth, we have omitted the reset procedure of the CUSUM. Note that the bound for b is tight, small deviations from \bar{b} lead to (boundedness) unboundedness of S_k . Next, for desired false alarm rates $\mathcal{A}^* \in \{0.25, 0.10, 0.02\}$, we compute the corresponding thresholds $\tau = \tau^*$ using Theorem 2 and Remark 4. For these thresholds, in Table 1, we present the actual false alarm rate \mathcal{A} (obtained by simulation) and the desired \mathcal{A}^* . Note that the difference between \mathcal{A} and \mathcal{A}^* is less than 0.05 in all cases.

In Figure 3, we present the evolution of $\|e_k^\delta\|$ when both the chi-squared and the CUSUM are deployed for attack detection and the attack sequences δ_k are zero-alarm attacks of the form introduced in (30) and (36), respectively. We consider two attack sequences, first, $\bar{\delta}_k^\alpha = \bar{\delta}_k^{\alpha,1} = \sqrt{\alpha/m\bar{\delta}_1}$, $\bar{\delta}_k^b = \bar{\delta}_k^{b,1} = \sqrt{b/m\bar{\delta}_1}$, $\bar{\tau}_{k^*} = \bar{\tau}_{k^*}^1 = \sqrt{(\tau + b - S_{k-1})/m\bar{\delta}_1}$, and $\bar{\delta}_1 = \mathbf{1}_{m \times 1}$. The second attack is $\bar{\delta}_k^\alpha = \bar{\delta}_k^{\alpha,2} = \sqrt{\alpha\bar{\delta}_2}$, $\bar{\delta}_k^b = \bar{\delta}_k^{b,2} = \sqrt{b\bar{\delta}_2}$, and $\bar{\tau}_{k^*} = \bar{\tau}_{k^*}^2 = \sqrt{\tau + b - S_{k-1}\bar{\delta}_2}$, where $\bar{\delta}_2$ denotes

b/\bar{b}	$\mathcal{A}^* = 0.25$		$\mathcal{A}^* = 0.10$		$\mathcal{A}^* = 0.02$	
	$\tau = \tau^*$	\mathcal{A} (Simul.)	$\tau = \tau^*$	\mathcal{A} (Simul.)	$\tau = \tau^*$	\mathcal{A} (Simul.)
1.05	1.0282	0.2041	3.9602	0.0899	12.3208	0.0196
1.15	0.6872	0.2010	3.3699	0.0885	10.0327	0.0184
2.00	—	—	0.2528	0.0953	4.1002	0.0202

Table 1. Simulation Experiments. Results from Theorem 2 and Remark 1.

the unitary singular vector corresponding to the largest singular value of $(I - F)^{-1}L\Sigma^{\frac{1}{2}}$. It can be proved that this selection of $\bar{\delta}_2$ maximizes the steady state value of $\|e_k^\delta\|$. For the CUSUM, we select $b = 2\bar{b} = 6$ and $\tau = \tau^* = 4.1002$ such that $\mathcal{A} \approx \mathcal{A}^* = 0.02$ (see Table 1). Likewise, we select $\alpha = \alpha^* = 2P^{-1}(\frac{2}{\tau}, 1 - 0.02) = 9.83$ such that, according to Theorem 3, $\mathcal{A} = \mathcal{A}^* = 0.02$. The attacks are induced at $k = k^* = 2 \times 10^3$. Note that, as stated in Proposition 1 and Proposition 2, given that $\rho[F] < 1$, inequalities (34) and (40) are satisfied for both attacks. Moreover, as mentioned in Section V-B, we expect that the CUSUM leads to smaller steady state deviation on $\|e_k\|$ because $b < \alpha$ and $\bar{\delta}_k^{\alpha, i} = a\bar{\delta}_k^{b, i}$, $i = 1, 2$ for some $a \in \mathbb{R}_{>0}$. This is exactly what we see in Figure 3. Note that the ratio $\bar{\gamma}_{\chi^2}/\bar{\gamma}_{CS} = 1.28$ is fixed by our choice of false alarm rate and bias. In Figure 4, we depict the evolution of the ratio $\bar{\gamma}_{\chi^2}/\bar{\gamma}_{CS}$ versus the false alarm rate \mathcal{A} for different values of CUSUM bias b .

VII. CONCLUSIONS

In this paper, for a class of stochastic linear time-invariant systems, we have characterized a model-based CUSUM proce-

cedure for identifying compromised sensors. In particular, steady state Kalman filters have been proposed to estimate the state of the physical process; then, these estimates have been used to construct residual variables (between sensor measurements and estimations) which drive the CUSUM procedure. Using stability results for stochastic systems and Markov chain approximations of the CUSUM sequence, we have derived systematic tools for tuning the CUSUM procedure such that mean square boundedness of the CUSUM sequence is guaranteed and the desired false alarm rate is fulfilled. For a class of zero-alarm attacks, we have characterized the performance of the proposed CUSUM procedure in terms of the effect that the attack sequence can induce on the system dynamics. Then, we have compared this performance against the one obtained using chi-squared procedures. For the linearized model of the chemical reactor considered in [19], [18], by means of a simulation study, we have showed that our tools are useful for tuning the CUSUM procedure and provide accurate predictions about the performance of the detection scheme.

$$\left\{ \begin{array}{l} F = \begin{pmatrix} 0.8353 & 0 & 0 & 0 \\ 0 & 0.8324 & 0 & 0.0031 \\ 0 & 0.0001 & 0.1633 & 0 \\ 0 & 0.0280 & 0.0172 & 0.9320 \end{pmatrix} G = \begin{pmatrix} 0.0458 & 0 & 0 \\ 0 & 0.0457 & 0 \\ 0 & 0 & 0.0231 \\ 0 & 0.0007 & 0.0006 \end{pmatrix}, C = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \\ L = \begin{pmatrix} 0.8271 & 0 & 0 \\ 0 & 0.8243 & 0.0002 \\ 0 & 0.0002 & 0.1619 \\ 0 & 0.0481 & 0.0543 \end{pmatrix}, R_2 = 0.01 \times I_3, R_0 = R_1 = I_4, \Sigma = \begin{pmatrix} 1.0169 & 0 & 0 \\ 0 & 1.0169 & 0.0001 \\ 0 & 0.0001 & 1.0105 \end{pmatrix}. \end{array} \right. \quad (42)$$

REFERENCES

- [1] F. Pasqualetti, F. D'Àurfler, and F. Bullo, "Attack detection and identification in cyber-physical systems," *IEEE Transactions on Automatic Control*, vol. 58, pp. 2715–2729, 2013.
- [2] Y. Mo, E. Garone, A. Casavola, and B. Sinopoli, "False data injection attacks against state estimation in wireless sensor networks," in *Decision and Control (CDC), 2010 49th IEEE Conference on*, 2010, pp. 5967–5972.
- [3] C. Kwon, W. Liu, and I. Hwang, "Security analysis for cyber-physical systems against stealthy deception attacks," in *American Control Conference (ACC), 2013*, 2013, pp. 3344–3349.
- [4] F. Miao, Q. Zhu, M. Pajic, and G. J. Pappas, "Coding sensor outputs for injection attacks detection," in *Decision and Control (CDC), 2014 IEEE 53rd Annual Conference on*, 2014, pp. 5776–5781.
- [5] C.-Z. Bai and V. Gupta, "On kalman filtering in the presence of a compromised sensor: Fundamental performance bounds," in *American Control Conference (ACC), 2014*, 2014, pp. 3029–3034.
- [6] C. Z. Bai, F. Pasqualetti, and V. Gupta, "Security in stochastic control systems: Fundamental limitations and performance bounds," in *American Control Conference (ACC), 2015*, 2015, pp. 195–200.
- [7] A. Cárdenas, S. Amin, Z. Lin, Y. Huang, C. Huang, and S. Sastry, "Attacks against process control systems: Risk assessment, detection, and response," in *Proceedings of the 6th ACM Symposium on Information, Computer and Communications Security*, 2011, pp. 355–366.
- [8] C. Murguía and J. Ruths, "Cusum and chi-squared attack detection of compromised sensors," in *proceedings of the IEEE Multi-Conference on Systems and Control (MSC)*, 2016.
- [9] —, "Characterization of a cusum model-based sensor attack detector," in *proceedings of the 55th IEEE Conference Decision and Control (CDC)*, 2016.
- [10] C. Murguía, N. van de Wouw, and J. Ruths, "Cusum and chi-squared attack detection of compromised sensors," in *proceedings of the 20th World Congress of the International Federation of Automatic Control (IFAC)*, accepted, 2017.
- [11] M. Ross, *Introduction to Probability Models, Ninth Edition*. Orlando, FL, USA: Academic Press, Inc., 2006.
- [12] H. Nijmeijer and A. van der Schaft, *Nonlinear dynamical control systems*. New York: Springer, 1990.
- [13] Z. Guo, D. Shi, K. H. Johansson, and L. Shi, "Optimal linear cyber-attack on remote state estimation," *IEEE Transactions on Control of Network Systems*, vol. 4, no. 1, pp. 4–13, 2017.

- [14] F. Gustafsson, *Adaptive Filtering and Change Detection*. West Sussex, Chichester, England: John Wiley and Sons, LTD, 2000.
- [15] K. Cheolhyeon, Y. Scott, and H. Inseok, "Real-time safety assessment of unmanned aircraft systems against stealthy cyber attacks," *Journal of Aerospace Information Systems*, vol. 13, pp. 27–45, 2015.
- [16] D. I. Urbina, J. A. Giraldo, A. A. Cardenas, N. O. Tippenhauer, J. Valente, M. Faisal, J. Ruths, R. Candell, and H. Sandberg, "Limiting the impact of stealthy attacks on industrial control systems," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '16, 2016, pp. 1092–1105.
- [17] D. I. Urbina, J. Giraldo, A. A. Cardenas, J. Valente, M. Faisal, N. O. Tippenhauer, J. Ruths, and H. Sandberg, "Survey and new directions for physics-based attack detection in control systems," in *NIST Standard Reference Simulation Website, NIST Standard Reference Database Number 16-010, National Institute of Standards and Technology, Gaithersburg MD, 20899*, 2016.
- [18] J. Chen and R. Patton, *Robust Model-Based Fault Diagnosis for Dynamic Systems*. Springer Publishing Company, Incorporated, 2012.
- [19] K. Watanabe and D. M. Himmelblau, "Fault diagnosis in nonlinear chemical processes. part ii. application to a chemical reactor," *AICHE Journal*, vol. 29, 1983.
- [20] T. T. Ha, *Theory and Design of Digital Communication Systems*. Cambridge University Press, 2010.
- [21] L. Ljung and T. Glad, *Modeling of Dynamic Systems*. Englewood Cliffs: PTR Prentice Hall, 1994.
- [22] K. J. Aström and B. Wittenmark, *Computer-controlled Systems (3rd Ed.)*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1997.
- [23] B. Anderson and J. Moore, *Optimal Filtering*. Englewood Cliffs, NJ: Prentice-Hall, 1979.
- [24] A. Wald, "Sequential tests of statistical hypotheses," *Ann. Math. Statist.*, vol. 16, pp. 117–186, 1945.
- [25] A. Willsky, "A survey of design methods for failure detection in dynamic systems," *Automatica*, vol. 12, pp. 601 – 611, 1976.
- [26] E. Page, "Continuous inspection schemes," *Biometrika*, vol. 41, pp. 100–115, 1954.
- [27] M. Basseville, "Detecting changes in signals and systems - a survey," *Automatica*, vol. 24, pp. 309 – 326, 1988.
- [28] J. Gertler, "Survey of model-based failure detection and isolation in complex plants," *Control Systems Magazine, IEEE*, vol. 8, pp. 3–11, 1988.
- [29] R. Agniel and E. Jury, "Almost sure boundedness of randomly sampled systems," *SIAM Journal on Control*, vol. 9, pp. 372–384, 1971.
- [30] T. Tarn and R. Yona, "Observers for nonlinear stochastic systems," *Automatic Control, IEEE Transactions on*, vol. 21, pp. 441–448, 1976.
- [31] C. van Dobben de Bruyn, *Cumulative sum tests : theory and practice*. London : Griffin, 1968.
- [32] B. Adams, W. Woodall, and C. Lowry, "The use (and misuse) of false alarm probabilities in control chart design," *Frontiers in Statistical Quality Control 4*, pp. 155–168, 1992.
- [33] R. A. Khan, "Wald's approximations to the average run length in cusum procedures," *Journal of Statistical Planning and Inference*, vol. 2, pp. 63 – 77, 1978.
- [34] C. Park, "A corrected wiener process approximation for cusum arls," *Sequential Analysis*, vol. 6, pp. 257–265, 1987.
- [35] M. Reynolds, "Approximations to the average run length in cumulative sum control charts," *Technometrics*, vol. 17, pp. 65–71, 1975.
- [36] C. Champ and S. Rigdon, "A comparison of the markov chain and the integral equation approaches for evaluating the run length distribution of quality control charts," *Communications in Statistics - Simulation and Computation*, vol. 20, pp. 191–204, 1991.
- [37] D. A. E. D. Brook, "An approach to the probability distribution of cusum run length," *Biometrika*, vol. 59, no. 3, pp. 539–549, 1972.
- [38] A. Luceno and J. Puig-vey, "Evaluation of the run-length probability distribution for cusum charts: Assessing chart performance," *Technometrics*, vol. 42, pp. 411–416, 2000.
- [39] W. Woodall, "The distribution of the run length of one-sided cusum procedures for continuous random variables," *Technometrics*, vol. 25, pp. 295–301, 1983.
- [40] S. Meyn and R. Tweedie, *Markov Chains and Stochastic Stability*. Springer-Verlag, 1993.
- [41] R. A. Horn and C. R. Johnson, *Matrix Analysis*, 2nd ed. New York, NY, USA: Cambridge University Press, 2012.
- [42] T. Dayar, "On moments of discrete phase-type distributions," in *Formal Techniques for Computer Systems and Business Processes*, M. Bravetti, L. Kloul, and G. Zavattaro, Eds. Versailles, France: Springer, 2005, pp. 51–63.
- [43] "On characterizations of the input-to-state stability property," *Systems and Control Letters*, vol. 24, pp. 351 – 359, 1995.
- [44] "Input-to-state stability for discrete-time nonlinear systems," *Automatica*, vol. 37, pp. 857 – 869, 2001.
- [45] G. Stewart, *Matrix Algorithms*. Society for Industrial and Applied Mathematics, 2001.

APPENDIX A PROOF OF THEOREM 1

Define the functions $V_k := S_k^2$ and $\Delta V_k := E_{S_k}[V_{k+1}] - V_k$. Along (18) with distance measure z_k and independent of τ , we have that

$$\Delta V_k = E_{S_k}[(S_k + z_{k+1} - b)^2] - S_k^2, \quad (43)$$

where $\rho^{+2} := \max(0, \rho)^2$. Note that, by construction, $S_k, z_k \in \mathbb{R}_{\geq 0}$ for all $k \in \mathbb{N}$. First, consider $S_k \in [b, \infty)$, it follows that

$$(S_k + z - b)^2 = (S_k + z - b)^2,$$

for all $z \in \mathbb{R}_{\geq 0}$; therefore

$$E_{S_k}[V_{k+1}] = E_{S_k}[(S_k + z_{k+1} - b)^2], \quad (44)$$

for $S_k \in [b, \infty)$ and $z_{k+1} \in \mathbb{R}_{\geq 0}$. Next, consider $S_k \in [0, b)$ which implies $S_k - b < 0$, then

$$z \in [0, b - S_k] \Rightarrow (S_k + z - b)^2 = 0,$$

$$z \in (b - S_k, \infty) \Rightarrow (S_k + z - b)^2 = (S_k + z - b)^2.$$

It follows that

$$E_{S_k}[V_{k+1}] \leq E_{S_k}[(S_k + z_{k+1} - b)^2], \quad (45)$$

for $S_k \in [0, b)$ and $z_{k+1} \in \mathbb{R}_{\geq 0}$. Using (44), (45), and independence between S_k and z_{k+1} , we have that

$$\begin{aligned} E_{S_k}[V_{k+1}] &\leq E_{S_k}[(S_k + z_{k+1} - b)^2] \\ &= (S_k - b)^2 + 2(S_k - b)E[z_{k+1}] + E[z_{k+1}^2], \end{aligned} \quad (46)$$

for $S_k, z_{k+1} \in \mathbb{R}_{\geq 0}$. Using (16) and the relation:

$$\text{var}[z_k] = E[z_k^2] - E[z_k]^2,$$

inequality (46) amounts to

$$E_{S_k}[V_{k+1}] \leq (S_k - b)^2 + 2(S_k - b + 1)m + m^2, \quad (47)$$

and, therefore,

$$\begin{aligned} \Delta V_k &= E_{S_k}[V_{k+1}] - S_k^2 \\ &\leq -2(b - m)S_k + (b - m)^2 + 2m, \end{aligned} \quad (48)$$

for $S_k, z_{k+1} \in \mathbb{R}_{\geq 0}$. From (48), given that $b \in \mathbb{R}_{> 0}$ and $S_k \in \mathbb{R}_{> 0}$ by construction, it is easy to verify that

$$\Delta V_k < 0 \Leftrightarrow b > \bar{b} := m \text{ and } S_k > \bar{S}, \quad \bar{S} := \frac{(b-m)^2 + 2m}{2(b-m)}. \quad (49)$$

Therefore, $b \in (\bar{b}, \infty)$ implies:

$$\begin{cases} \Delta V_k < 0 & \text{for } S_k \in (\bar{S}, \infty), \\ \Delta V_k \geq 0 & \text{for } S_k \in [0, \bar{S}]. \end{cases} \quad (50)$$

Recall that $\Delta V_k := E_{S_k}[V_{k+1}] - V_k$. Assume that for some $k = k^* \in \mathbb{N}$, $S_{k^*} \in (\bar{S}, \infty)$; then, from (50), $\Delta V_{k^*} < 0$ and consequently

$$E_{S_{k^*}}[V_{k^*+1}] < V_{k^*}. \quad (51)$$

Next, for $k = k^* + 1$, let $S_{k^*+1} \in (\bar{S}, \infty)$, then

$$E_{S_{k^*+1}}[V_{k^*+2}] < V_{k^*+1}. \quad (52)$$

Using (51), (52), and the property

$$E_{S_{k^*}}[V_{k^*+2}] = E_{S_{k^*}}[E_{S_{k^*+1}}[V_{k^*+2}]], \quad (53)$$

we have

$$E_{S_{k^*}}[V_{k^*+2}] < E_{S_{k^*}}[V_{k^*+1}] < V_{k^*}. \quad (54)$$

Continuing this way, we obtain

$$V_{k^*} > E_{S_{k^*}}[V_{k^*+1}] > \dots > E_{S_{k^*}}[V_{k^*+n}], \quad (55)$$

for $n \geq 2$ and $k = \{k^*, k^* + 1, \dots, k^* + n\}$ such that $S_k \in (\bar{S}, \infty)$. Therefore, from (55), it can be concluded that the second moment $E_{S_{k^*}}[V_k] = E_{S_{k^*}}[S_k^2]$ decreases monotonically for $S_k \in (\bar{S}, \infty)$ and $b \in (\bar{b}, \infty)$; and consequently, $E_{S_{k^*}}[S_k^2] < \infty$. Next, assume that for some $k = k^* \in \mathbb{N}$, $S_{k^*} \in [0, \bar{S}]$ and $b \in (\bar{b}, \infty)$; then, from (50), $\Delta V_{k^*} \geq 0$ and, by (48), it follows that

$$\Delta V_{k^*} \leq -2(b-m)S_{k^*} + (b-m)^2 + 2m \geq 0. \quad (56)$$

Since $S_{k^*} \in [0, \bar{S}]$, there always exist constants $a \in (0, 1)$ and $\beta \in \mathbb{R}_{>0}$ satisfying

$$\begin{aligned} \Delta V_{k^*} &\leq -2(b-m)S_{k^*} + (b-m)^2 + 2m \\ &\leq -aV_{k^*} + \beta, \quad \text{for } S_{k^*} \in [0, \bar{S}]. \end{aligned} \quad (57)$$

Given that $\Delta V_k = E_{S_k}[V_{k+1}] - V_k$, by (57), we have

$$E_{S_{k^*}}[V_{k^*+1}] \leq (1-a)V_{k^*} + \beta, \quad \text{for } S_{k^*} \in [0, \bar{S}]. \quad (58)$$

Next, for $k = k^* + 1$, let $S_{k^*+1} \in [0, \bar{S}]$, then

$$E_{S_{k^*+1}}[V_{k^*+2}] \leq (1-a)V_{k^*+1} + \beta. \quad (59)$$

By (58), (59), and the property

$$E_{S_{k^*}}[V_{k^*+2}] = E_{S_{k^*}}[E_{S_{k^*+1}}[V_{k^*+2}]], \quad (60)$$

we have

$$\begin{aligned} E_{S_{k^*}}[V_{k^*+2}] &\leq E_{S_{k^*}}[(1-a)V_{k^*+1} + \beta] \\ &= (1-a)E_{S_{k^*}}[V_{k^*+1}] + \beta \\ &\leq (1-a)^2V_{k^*} + (1-a)\beta + \beta. \end{aligned} \quad (61)$$

Continuing this way, we obtain

$$E_{S_{k^*}}[V_{k^*+n}] \leq (1-a)^nV_{k^*} + \beta \sum_{i=0}^{n-1} (1-a)^i, \quad (62)$$

for $n \geq 1$ and $k = \{k^*, k^* + 1, \dots, k^* + n\}$ such that $S_k \in [0, \bar{S}]$. Given that $a \in (0, 1)$ and $\sum_{i=0}^{k-1} (1-a)^i \leq \sum_{i=0}^{\infty} (1-a)^i = \frac{1}{a}$, then, as $n \rightarrow \infty$, it is satisfied that

$$E_{S_{k^*}}[S_{\infty}^2] \leq \frac{\beta}{a}. \quad (63)$$

Therefore, by (62) and (63), the second moment $E_{S_{k^*}}[S_k^2]$ does not grow unbounded for $S_k \in [0, \bar{S}]$ and $b \in (\bar{b}, \infty)$, i.e., $E_{S_{k^*}}[S_k^2] < \infty$. So far, combining the preliminary results presented above, we have proved that for all $S_k \in [0, b] \cup (b, \infty)$, the second moment is finite and either decreasing or uniformly bounded in $k \in \mathbb{N}$ provided that the conditions of Theorem 1 are satisfied. Hence, for the residual sequence $r_k \sim \mathcal{N}(\mathbf{0}, \Sigma)$, $b > \bar{b} \rightarrow E_{S_1}[S_k^2] < \infty$ for all $k \in \mathbb{N}$. ■