

# Security in Stochastic Control Systems: Fundamental Limitations and Performance Bounds

Cheng-Zong Bai, Fabio Pasqualetti, and Vijay Gupta

**Abstract**—This work proposes a novel metric to characterize the resilience of stochastic cyber-physical systems to attacks and faults. We consider a single-input single-output plant regulated by a control law based on the estimate of a Kalman filter. We allow for the presence of an attacker able to hijack and replace the control signal. The objective of the attacker is to maximize the estimation error of the Kalman filter – which in turn quantifies the degradation of the control performance – by tampering with the control input, while remaining undetected. We introduce a notion of  $\epsilon$ -stealthiness to quantify the difficulty to detect an attack when an arbitrary detection algorithm is implemented by the controller. For a desired value of  $\epsilon$ -stealthiness, we quantify the largest estimation error that an attacker can induce, and we analytically characterize an optimal attack strategy. Because our bounds are independent of the detection mechanism implemented by the controller, our information-theoretic analysis characterizes fundamental security limitations of stochastic cyber-physical systems.

## I. INTRODUCTION

Cyber-physical systems offer a variety of attack surfaces arising from the interconnection of different technologies and components. Depending on their resources and capabilities, attackers generally aim to deteriorate the functionality of the system, while avoiding detection for as long as possible [1].

Security of cyber-physical systems is a growing research area where, recently, different attack strategies and defense mechanisms have been characterized. While simple attacks have a straightforward implementation and impact, such as jamming control and communication channels [2], sophisticated ones may degrade the functionality of a system more severely [3], [4], and are more difficult to mitigate. In this work we measure the severity of attacks based on their effect on the control performance and on their level of stealthiness, that is, the difficulty of being detected from measurements. Intuitively, there exists a trade-off between the degradation of control performance and the level of stealthiness of an attack. Although this trade-off has previously been identified for specific systems and detection mechanisms [5], [6], [7], [8], a thorough analysis of the resilience of stochastic control systems to arbitrary attacks is still missing.

**Related works** For deterministic cyber-physical systems the concept of stealthiness of an attack is closely related to the control-theoretic notion of zero dynamics [9]. In particular, an attack is undetectable if and only if it excites only the zero

dynamics of an appropriately defined input-output system describing the system dynamics, the measurements available to a security monitors, and the variables compromised by the attacker [10], [11]. Thus, the question of stealthiness of an attack has a binary answer in deterministic systems. For stochastic cyber-physical systems, instead, the presence of process and measurements noise offers a smart attacker the additional possibility to tamper with sensor measurements and control inputs within the acceptable uncertainty levels, thereby making the detection task arbitrarily difficult.

Detectability of attacks in stochastic systems has received only initial attention from the research community, and there seem to be no agreement on an appropriate notion of stealthiness. Most works in this area consider detectability of attacks with respect to specific detection schemes, such as the classic bad data detection algorithm [12]. In our previous work [13], we proposed the notion of  $\epsilon$ -marginal stealthiness to quantify the stealthiness level with respect to the class of ergodic detectors. With respect to [13], in this work (i) we introduce a novel notion of stealthiness, namely  $\epsilon$ -stealthiness, that is independent of the attack detection algorithm and thus provides a fundamental measure of the stealthiness of attacks in stochastic control systems, and (ii) we explicitly characterize detectability and performance of  $\epsilon$ -stealthy attacks.

**Contributions** The contributions of this paper are threefold. First, we propose the notion of  $\epsilon$ -stealthiness to quantify detectability of attacks in stochastic cyber-physical systems. Our metric is motivated by the Chernoff-Stein Lemma in detection and information theories [14], and is universal, in the sense that it is independent of any specific detection mechanism employed by the controller. Second, we provide an achievable bound for the degradation of the minimum-mean-square estimation error caused by an  $\epsilon$ -stealthy attack, as a function of the system parameters, noise statistics, and information available to the attacker. Third and finally, we provide a closed-form expression of optimal  $\epsilon$ -stealthy attacks achieving the maximal degradation of the estimation error. These results characterize the trade-off between performance degradation that an attacker can induce, versus the fundamental limit of the detectability of the attack.

We focus on single-input single-output systems with an observer-based controller. However, our methods are general, and applicable to multiple-input multiple-output systems via a more involved technical analysis.

**Paper organization** Section II contains our mathematical formulation of the problem and our model of attacker. In Section III we discuss our metric to quantify the stealthiness

This material is based upon work supported in part by awards NSF ECCS-1405330 and ONR N00014-14-1-0816. Cheng-Zong Bai and Vijay Gupta are with the Department of Electrical Engineering, University of Notre Dame, IN 46556, {cbai, vgupta2}@nd.edu. Fabio Pasqualetti is with the Department of Mechanical Engineering, University of California, Riverside, CA 92521, fabiopas@engr.ucr.edu.

level of an attack. The main results of this paper are presented in Section IV, including a characterization of the largest perturbation caused by an  $\epsilon$ -stealthy attack, and a closed-form expression of optimal  $\epsilon$ -stealthy attacks. Section V contains our illustrative examples and numerical results. Finally, Section VI concludes the paper.

## II. SYSTEM AND ATTACK MODELS

In this section we detail our system and attack models. Throughout the paper, we let  $x_i^j$  denote the sequence  $\{x_n\}_{n=i}^j$ , and  $x \sim \mathcal{N}(\mu, \sigma^2)$  a Gaussian random variable with mean  $\mu$  and variance  $\sigma^2$ .

### A. System model

We consider the single-input single-output time-invariant system described by

$$x_{k+1} = ax_k + u_k + w_k, \quad y_k = cx_k + v_k, \quad (1)$$

where  $a, c \in \mathbb{R}$ ,  $c \neq 0$ ,  $w_1^\infty$  and  $v_1^\infty$  are random sequences representing process and measurement noise, respectively. We assume the sequences  $w_1^\infty$  and  $v_1^\infty$  to be independent and identically distributed (i.i.d.) Gaussian processes with  $w_k \sim \mathcal{N}(0, \sigma_w^2)$ ,  $v_k \sim \mathcal{N}(0, \sigma_v^2)$  for all  $k > 0$ . The control input  $u_k$  is generated based on a causal observer-based control policy, that is,  $u_k$  is a function of the measurement sequence  $y_1^k$ . In particular, the controller employs a Kalman filter [15], [16] to compute the Minimum-Mean-Squared-Error (MMSE) estimate  $\hat{x}_{k+1}$  of  $x_{k+1}$  from the measurements  $y_1^k$ . The Kalman filter reads as

$$\hat{x}_{k+1} = a\hat{x}_k + K_k(y_k - c\hat{x}_k) + u_k, \quad (2)$$

where the Kalman gain  $K_k$  and the mean squared error  $P_{k+1} \triangleq \mathbb{E}[(\hat{x}_{k+1} - x_{k+1})^2]$  can be calculated by the recursions

$$K_k = \frac{acP_k}{c^2P_k + \sigma_v^2}, \quad P_{k+1} = a^2P_k + \sigma_w^2 - \frac{a^2c^2P_k^2}{c^2P_k + \sigma_v^2}.$$

with the initial condition  $\hat{x}_1 = \mathbb{E}[x_1] = 0$  and  $P_1 = \mathbb{E}[x_1^2]$ . If the system (1) is detectable (i.e.,  $|a| < 1$  or  $c \neq 0$ ), then the Kalman filter converges to the steady state in the sense that  $\lim_{k \rightarrow \infty} P_k = P$  exists [16], where  $P$  can be obtained uniquely through the algebraic Riccati equation. For the ease of presentation, we assume that  $P_1 = P$ . Hence, we obtain a steady state Kalman filter with Kalman gain  $K_k = K$  and  $P_k = P$  at every time step  $k$ . The sequence  $z_1^\infty$  calculated as  $z_k \triangleq y_k - c\hat{x}_k$  is called the innovation sequence. Since we consider steady state Kalman filtering, the innovation sequence is an i.i.d. Gaussian process with  $z_k \sim \mathcal{N}(0, c^2P + \sigma_v^2)$ .

### B. Attack model

We consider an attacker capable of hijacking and replacing the control input  $u_1^\infty$  with an arbitrary signal  $\tilde{u}_1^\infty$ . Assume that the attacker knows the system parameters  $a, c, \sigma_w^2, \sigma_v^2$ . Let  $\mathcal{I}_k$  denote the information available to the attacker at time  $k$ . The attack input  $\tilde{u}_1^\infty$  is constructed based on

the system parameters and the attacker *information pattern*, which satisfies the following assumptions:

- (A1) the attacker knows the control input  $u_k$ , that is,  $u_k \in \mathcal{I}_k$  at all times  $k$ ;
- (A2) the information available to the attacker is non-decreasing, that is,  $\mathcal{I}_k \subseteq \mathcal{I}_{k+1}$ ;
- (A3)  $\mathcal{I}_k$  is independent of the  $w_k^\infty$  and  $v_{k+1}^\infty$  due to causality.

Attack scenarios satisfying assumptions (A1)–(A3) include:

- (i) the attacker knows the control input, that is,  $\mathcal{I}_k = \{u_1^k\}$ ;
- (ii) the attacker knows the control input and the state, that is,  $\mathcal{I}_k = \{u_1^k, x_1^k\}$ ;
- (iii) the attacker knows the control input and the (delayed) measurements received by the controller, that is,  $\mathcal{I}_k = \{u_1^k, \tilde{y}_1^{k-d}\}$  with  $d \leq 0$ ;
- (iv) the attacker knows the control input and take additional measurements  $\bar{y}_k$ , that is,  $\mathcal{I}_k = \{u_1^k, \bar{y}_1^k\}$ .

Let  $\tilde{y}_1^\infty$  be the sequence of measurements received by the controller in the presence of the attack  $\tilde{u}_1^\infty$ . Then,  $\tilde{y}_1^\infty$  is generated by the dynamics

$$x_{k+1} = ax_k + \tilde{u}_k + w_k, \quad \tilde{y}_k = cx_k + v_k. \quad (3)$$

Notice that, because the controller is unaware of the attack, the corrupted measurements  $\tilde{y}_1^\infty$ , and hence the attack input  $\tilde{u}_1^\infty$ , drive the Kalman filter (2) as an external input. Let  $\hat{x}_1^\infty$  be the estimate of the Kalman filter (2) in the presence of the attack  $\tilde{u}_1^\infty$ , which is obtained from the recursion

$$\hat{x}_{k+1} = a\hat{x}_k + K\tilde{z}_k + u_k,$$

with innovation is  $\tilde{z}_k \triangleq \tilde{y}_k - c\hat{x}_k$ . Notice that (i) the estimate  $\hat{x}_{k+1}$  is sub-optimal, because it is obtained by assuming the nominal control input, whereas the system is driven by the attack input, and (ii) the random sequence  $\tilde{z}_1^\infty$  need neither be stationary, nor zero mean, white or Gaussian, because the attack input is arbitrary.

Let  $\tilde{P}_{k+1} = \mathbb{E}[(\hat{x}_{k+1} - x_{k+1})^2]$  be the second moment of the estimation error  $\hat{x}_{k+1} - x_{k+1}$ , and assume that the attacker aims to maximize  $\tilde{P}_{k+1}$ . We consider the asymptotic behavior of  $\tilde{P}_{k+1}$  to measure the performance degradation induced by the attacker. Since the attack sequence is arbitrary, the sequence  $\tilde{P}_1^\infty$  may diverge. Accordingly, we consider the limit superior of arithmetic mean of the sequence  $\tilde{P}_1^\infty$  as given by

$$\tilde{P} \triangleq \limsup_{k \rightarrow \infty} \frac{1}{k} \sum_{n=1}^k \tilde{P}_n.$$

Notice that if the sequence  $\tilde{P}_1^\infty$  is convergent, then  $\lim_{k \rightarrow \infty} \tilde{P}_{k+1} = \tilde{P}$ , which equals the Cesàro mean<sup>1</sup> [14].

<sup>1</sup>The steady state assumption is made in order to obtain an i.i.d. innovation sequence. If the Kalman filter starts from an arbitrary initial condition  $P_1$ , then the innovation sequence is an independent, asymptotically identically distributed, Gaussian process. This identity guarantees that the results for the case of non-steady state Kalman filter coincide with the main results (i.e., Theorem 1 and Theorem 2) in this paper.

### III. ATTACK STEALTHINESS FOR STOCHASTIC SYSTEMS

In this section we motivate and define our notion of  $\epsilon$ -stealthiness of attacks. Notice that the system (3) with  $\sigma_w^2 = 0$  and  $\sigma_v^2 = 0$  (i.e., deterministic single-input single-output system) features no zero dynamics. Hence, every attack would be detectable [10]. However, the stochastic nature of the system provides an additional degree of freedom to the attacker, because the process noise and the measurement noise induce some uncertainty in the measurements. Building on this idea, we now formally define attack stealthiness. Consider the problem of detecting an attack from measurements. Notice that the detector must rely on the statistical properties of the received measurement sequence as compared with their expected model in (1). This can be formulated by the following binary hypothesis testing problem:

$$\begin{aligned} H_0 &: \text{No attack is in progress (the controller receives } y_1^k); \\ H_1 &: \text{Attack is in progress (the controller receives } \tilde{y}_1^k). \end{aligned}$$

Suppose that a detector is employed by the controller. Let  $p_k^F$  be the probability of false alarm (decide  $H_1$  when  $H_0$  is true) at time  $k$  and let  $p_k^D$  be the probability of detection (decide  $H_1$  when  $H_1$  is true) at time  $k$ . In detection theory, the performance of the detector can be characterized by the trade-off between  $p_k^F$  and  $p_k^D$ , namely, the Receiver Operating Characteristic (ROC) [17]. From the ROC perspective, the attack that is hardest to detect is the one for which, at every time  $k$ , there exists no detector that performs better than a random guess (e.g., to make a decision by flipping a coin) independent of the hypothesis. If a detector makes a decision via a random guess independent of the hypothesis, then the operating point of the ROC satisfies  $p_k^F = p_k^D$ .

**Definition 1: (Strict stealthiness)** The attack  $\tilde{u}_1^\infty$  is strictly stealthy if there exists no detector such that  $p_k^F < p_k^D$  at any  $k > 0$ .  $\square$

The reader may argue that strict stealthiness is a too restrictive notion of stealthiness for an attacker, and it significantly limits the set of stealthy attacks. In fact, the attacker may be satisfied with attack inputs that are difficult to detect, in the sense that the detector would need to collect more measurements to make a decision with a desired operating point of ROC. Although it is impractical to compute the exact values of these two probabilities for an arbitrary detector at every time  $k$ , we are able to apply the techniques in detection theory and information theory to obtain bounds for  $p_k^F$  and  $p_k^D$ . A classical example is the Chernoff-Stein Lemma [14]. This lemma characterizes the asymptotic exponent of  $p_k^F$ , while  $p_k^D$  can be arbitrary. Motivated by Chernoff-Stein Lemma, we propose the following notion of  $\epsilon$ -stealthiness.

**Definition 2: ( $\epsilon$ -stealthiness)** Let  $\epsilon > 0$  and  $0 < \delta < 1$ . The attack  $\tilde{u}_1^\infty$  is  $\epsilon$ -stealthy if there exists no detector such that the following two conditions can be satisfied simultaneously:

- (i) The detector operates with  $0 < 1 - p_k^D \leq \delta$  at all times  $k$ .
- (ii) The probability of false alarm  $p_k^F$  converges to zero exponentially fast with rate greater than  $\epsilon$  as  $k$  grows.

In other words, for any detector that satisfies  $0 < 1 - p_k^D \leq \delta$  for all times  $k$ , it holds

$$\limsup_{k \rightarrow \infty} -\frac{1}{k} \log p_k^F \leq \epsilon. \quad (4)$$

$\square$

Definition 2 provides a characterization of the detectability for  $\epsilon$ -stealthy attacks. We now provide a sufficient condition and a necessary condition for an attack to be  $\epsilon$ -stealthy, which rely on the Kullback-Leibler divergence (or relative entropy) [14], [18] defined as follows.

**Definition 3: (Kullback-Leibler divergence)** Let  $x_1^k$  and  $y_1^k$  be two random sequences with joint probability density functions  $f_{x_1^k}$  and  $f_{y_1^k}$ , respectively. The Kullback-Leibler Divergence (KLD) between  $x_1^k$  and  $y_1^k$  equals

$$D(x_1^k \| y_1^k) = \int_{-\infty}^{\infty} \log \frac{f_{x_1^k}(\xi_1^k)}{f_{y_1^k}(\xi_1^k)} f_{x_1^k}(\xi_1^k) d\xi_1^k. \quad (5)$$

$\square$

The KLD is a non-negative measure of the dissimilarity between two probability density functions. It should be observed that  $D(x_1^k \| y_1^k) = 0$  if  $f_{x_1^k} = f_{y_1^k}$ . Also, the KLD is generally not symmetric, that is,  $D(x_1^k \| y_1^k) \neq D(y_1^k \| x_1^k)$ . Using the Chernoff-Stein Lemma, we can provide a sufficient condition for an attack to be  $\epsilon$ -stealthy.

**Lemma 1: (Sufficient condition for  $\epsilon$ -stealthiness)** Let  $\tilde{y}_1^\infty$  be the random sequence generated by the attack  $\tilde{u}_1^\infty$ . Let  $\tilde{y}_1^\infty$  be ergodic and satisfy

$$\lim_{k \rightarrow \infty} \frac{1}{k} D(\tilde{y}_1^k \| y_1^k) \leq \epsilon. \quad (6)$$

Then, the attack  $\tilde{u}_1^\infty$  is  $\epsilon$ -stealthy.

*Proof:* We apply the Chernoff-Stein Lemma for ergodic measurements (e.g., see [19]). For such an attack  $\tilde{u}_1^\infty$ , given  $0 < 1 - p_k^D \leq \delta$  where  $0 < \delta < 1$ , the best achievable exponent of  $p_k^F$  is given by  $\lim_{k \rightarrow \infty} \frac{1}{k} D(\tilde{y}_1^k \| y_1^k)$ . For any detector, we obtain

$$\limsup_{k \rightarrow \infty} -\frac{1}{k} \log p_k^F \leq \lim_{k \rightarrow \infty} \frac{1}{k} D(\tilde{y}_1^k \| y_1^k) \leq \epsilon.$$

By Definition 2, the attack is  $\epsilon$ -stealthy.  $\blacksquare$

Next, we provide a necessary condition for an attack to be  $\epsilon$ -stealthy.

**Lemma 2: (Necessary condition for  $\epsilon$ -stealthiness)** Let the attack  $\tilde{u}_1^\infty$  be  $\epsilon$ -stealthy. Then

$$\limsup_{k \rightarrow \infty} \frac{1}{k} D(\tilde{y}_1^k \| y_1^k) \leq \epsilon. \quad (7)$$

*Proof:* The proof can be found in [20].  $\blacksquare$

We conclude this section with a method to compute the KLD between the sequences  $\tilde{y}_1^k$  and  $y_1^k$ . For observed-based controllers, note that  $z_k$  and  $\tilde{z}_k$  are invertible functions of  $y_1^k$  and  $\tilde{y}_1^k$ , respectively. Recall from the invariance properties of the KLD [18] that, for every  $k > 0$ ,

$$D(\tilde{y}_1^k \| y_1^k) = D(\tilde{z}_1^k \| z_1^k).$$

Moreover,  $z_1^\infty$  is an i.i.d. Gaussian random sequence with  $z_k \sim \mathcal{N}(0, \sigma_z^2)$ . From (5) we obtain

$$\frac{1}{k} D(\tilde{z}_1^k \| z_1^k) = -\frac{1}{k} h(\tilde{z}_1^k) + \frac{1}{2} \log(2\pi\sigma_z^2) + \frac{1}{k} \sum_{n=1}^k \frac{\mathbb{E}[\tilde{z}_n^2]}{2\sigma_z^2}, \quad (8)$$

where  $h(\tilde{z}_1^k) = \int_{-\infty}^{\infty} -f_{\tilde{z}_1^k}(\xi_1^k) \log f_{\tilde{z}_1^k}(\xi_1^k) d\xi_1^k$  is the differential entropy of  $\tilde{z}_1^k$  [14].

#### IV. PERFORMANCE BOUNDS AND LIMITATIONS

We are interested in the maximal performance degradation  $\tilde{P}$  that an  $\epsilon$ -stealthy attack may induce. We present such a fundamental limit in two parts: the converse statement that gives an upper bound for  $\tilde{P}$  as induced by an  $\epsilon$ -stealthy attack, and the achievability result that provides an attack that achieves the upper bound of the converse result.

**Theorem 1: (Converse)** Consider the system stated in (1). Let the sequence  $\mathcal{I}_1^\infty$  satisfy assumptions (A1)–(A3). Let  $\tilde{u}_1^\infty$  be an  $\epsilon$ -stealthy attack generated by  $\mathcal{I}_1^\infty$ . Then, the estimation error induced by the attacker satisfies

$$\tilde{P} = \limsup_{k \rightarrow \infty} \frac{1}{k} \sum_{n=1}^k \tilde{P}_n \leq \bar{\delta}(\epsilon)P + \frac{(\bar{\delta}(\epsilon) - 1)\sigma_v^2}{c^2} \quad (9)$$

where the function  $\bar{\delta}: [0, \infty) \rightarrow [1, \infty)$  is such that

$$\bar{\delta}(D) = 2D + 1 + \log \bar{\delta}(D). \quad (10)$$

*Proof:* Observe that  $\tilde{z}_k = \tilde{y}_k - c\hat{x}_k = c(x_k - \hat{x}_k) + v_k$ , and  $(x_k - \hat{x}_k)$  is independent of  $v_k$ . We have

$$\mathbb{E}[\tilde{z}_k^2] = c^2 \tilde{P}_k + \sigma_v^2. \quad (11)$$

Since  $\sigma_v^2$  is a constant and  $c^2 > 0$ , we can represent  $\tilde{P}$  in terms of  $\mathbb{E}[\tilde{z}_k^2]$ . From (8), we have

$$\begin{aligned} & \frac{1}{2} \cdot \frac{1}{k} \sum_{n=1}^k \frac{\mathbb{E}[\tilde{z}_n^2]}{\sigma_z^2} \\ &= \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) - \frac{1}{2} \log(2\pi\sigma_z^2) + \frac{1}{k} h(\tilde{z}_1^k) \\ &\leq \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) - \frac{1}{2} \log(2\pi\sigma_z^2) + \frac{1}{k} \sum_{n=1}^k h(\tilde{z}_n) \end{aligned} \quad (12)$$

$$\leq \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) - \frac{1}{2} \log(2\pi\sigma_z^2) + \frac{1}{k} \sum_{n=1}^k \frac{1}{2} \log(2\pi e \mathbb{E}[\tilde{z}_n^2]) \quad (13)$$

$$\begin{aligned} &= \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) + \frac{1}{2} + \frac{1}{2} \log \left( \prod_{n=1}^k \frac{\mathbb{E}[\tilde{z}_n^2]}{\sigma_z^2} \right)^{\frac{1}{k}} \\ &\leq \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) + \frac{1}{2} + \frac{1}{2} \log \left( \frac{1}{k} \sum_{n=1}^k \frac{\mathbb{E}[\tilde{z}_n^2]}{\sigma_z^2} \right), \end{aligned} \quad (14)$$

where the inequalities (12) is due to the subadditivity of differential entropy [14, Corollary 8.6.1], the inequality (13) is a consequence of the maximum entropy theorem [14, Theorem 8.6.5], and the inequality (14) follows from the

arithmetic mean and geometric mean inequality. Consider the following maximization problem

$$\begin{aligned} & \max_{x \in \mathbb{R}} \quad x, \\ & \text{subject to} \quad \frac{1}{2}x - D - \frac{1}{2} \leq \frac{1}{2} \log x, \end{aligned} \quad (15)$$

where  $D \geq 0$ . Since the logarithm function is concave, the feasible region of  $x$  in (15) is a closed interval upper bounded by  $\bar{\delta}(D)$  as defined in (10); see Fig. 1. Thus, the maximum in (15) is  $\bar{\delta}(D)$ . By (14) and the maximization problem (15), we obtain

$$\frac{1}{k} \sum_{n=1}^k \frac{\mathbb{E}[\tilde{z}_n^2]}{\sigma_z^2} \leq \bar{\delta} \left( \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) \right). \quad (16)$$

From (11) and (16) we obtain

$$\begin{aligned} \tilde{P} &= \limsup_{k \rightarrow \infty} \frac{1}{k} \sum_{n=1}^k \tilde{P}_n = \limsup_{k \rightarrow \infty} \frac{1}{k} \sum_{n=1}^k \frac{\mathbb{E}[\tilde{z}_n^2] - \sigma_v^2}{c^2} \\ &\leq \limsup_{k \rightarrow \infty} \frac{\bar{\delta} \left( \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) \right) \sigma_z^2 - \sigma_v^2}{c^2} \end{aligned} \quad (17)$$

$$= \frac{\bar{\delta} \left( \limsup_{k \rightarrow \infty} \frac{1}{k} D(\tilde{z}_1^k \| z_1^k) \right) \sigma_z^2 - \sigma_v^2}{c^2} \quad (18)$$

$$\leq \frac{\bar{\delta}(\epsilon) \sigma_z^2 - \sigma_v^2}{c^2}, \quad (19)$$

where the inequality (17) can be obtained by the definition of limit superior, the equality (18) is due to the continuity and monotonicity of the function  $\bar{\delta}$ , and the inequality (19) follows from Lemma 2. Finally, the desired result is obtained by substituting  $\sigma_z^2 = c^2 P + \sigma_v^2$  into (19). ■

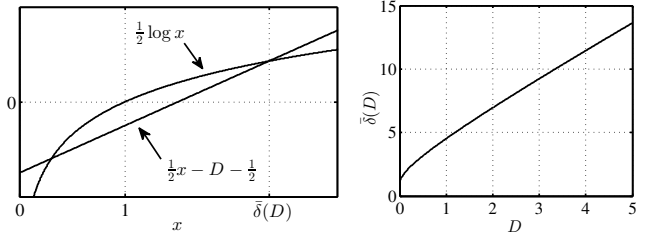


Fig. 1. Illustrations for the optimization problem (15) and the function  $\bar{\delta}: [0, \infty) \rightarrow [1, \infty)$  defined in (10). Notice that the function  $\bar{\delta}$  is continuous and monotonically increasing.

**Remark 1: (Effect of strictly stealthy attacks)** Strictly stealthy attacks do not degrade the performance of the Kalman filter. To see this, notice that if an attack is strictly stealthy then  $D(\tilde{y}_1^k \| y_1^k) = 0$  for all  $k > 0$  (this is a consequence of Definition 1 and the Neyman-Pearson Lemma [17]). Moreover, by using (11), (16), and the fact that  $\bar{\delta}(0) = 1$  whenever  $D(\tilde{z}_1^k \| z_1^k) = 0$  for all  $k > 0$ , we obtain  $\mathbb{E}[\tilde{z}_k^2] = c^2 \tilde{P}_k + \sigma_v^2 \leq c^2 P + \sigma_v^2$ . Consequently  $\tilde{P}_k \leq P$ , that is, the mean squared error of the Kalman filter under attack is less or equal to the minimum mean squared error in the absence of attacks. □

In the next theorem we construct an  $\epsilon$ -stealthy attack that achieves the upper bound in Theorem 1.

**Theorem 2: (Achievability)** Let  $\zeta_1^\infty$  be an i.i.d. sequence of random variables  $\zeta_k \sim \mathcal{N}(0, \frac{\sigma_z^2}{c^2}(\bar{\delta}(\epsilon) - 1))$  independent of  $\{x_1^k, \tilde{y}_1^k, \mathcal{I}_1^k\}$ , and let the attack be defined as

$$\tilde{u}_k = u_k - (a - Kc)\zeta_{k-1} + \zeta_k, \quad (20)$$

with  $\zeta_0 = 0$ . Then, the attack  $\tilde{u}_1^\infty$  is  $\epsilon$ -stealthy and it achieves the converse result in Theorem 1, that is,

$$\tilde{P} = \lim_{k \rightarrow \infty} \frac{1}{k} \sum_{n=1}^k \tilde{P}_n = \bar{\delta}(\epsilon)P + \frac{(\bar{\delta}(\epsilon) - 1)\sigma_v^2}{c^2},$$

where the function  $\bar{\delta} : [0, \infty) \rightarrow [1, \infty)$  satisfies (10).

*Proof:* For the ease of analysis and without affecting generality, we assume that the attack  $\tilde{u}_1^\infty$  is generated by an attacker with the information pattern  $\mathcal{I}_1^\infty$ , with  $\mathcal{I}_k = \{u_1^k, \tilde{y}_1^k\}$  for every  $k > 0$ .

We first show that the upper bound (9) is achieved by the attack. Notice that the attacker implements the Kalman filter  $\hat{x}_{k+1}^A = a\hat{x}_k^A + Kz_k^A + \tilde{u}_k$  with the initial condition  $\hat{x}_1^A = 0$  where  $z_k^A = \tilde{y}_k - c\hat{x}_k^A$ . Thus,  $\hat{x}_{k+1}^A$  is the MMSE estimate of the state with the mean squared error  $\mathbb{E}[(\hat{x}_{k+1}^A - x_{k+1})^2] = P$  when  $\mathcal{I}_k$  is given. Note that  $\tilde{z}_k$  can be expressed as

$$\tilde{z}_k = \tilde{y}_k - c\hat{x}_k = \tilde{y}_k - c\hat{x}_k^A + c(\hat{x}_k^A - \hat{x}_k) = z_k^A - c\tilde{e}_k, \quad (21)$$

where  $\tilde{e}_k = \hat{x}_k - \hat{x}_k^A$ . In addition, the dynamics of  $\tilde{e}_k$  are given by

$$\begin{aligned} \tilde{e}_{k+1} &= (a\hat{x}_k + K\tilde{z}_k + u_k) - (a\hat{x}_k^A + Kz_k^A + \tilde{u}_k) \\ &= (a - Kc)\tilde{e}_k + (a - Kc)\zeta_{k-1} - \zeta_k, \end{aligned} \quad (22)$$

and the initial condition is  $\tilde{e}_1 = 0$ . Equation (22) implies that  $\tilde{e}_{k+1} = -\zeta_k$  for every  $k > 0$ . Further, for every  $k > 0$ ,  $\tilde{P}_{k+1}$  can be expressed as

$$\begin{aligned} \tilde{P}_{k+1} &= \mathbb{E}[(\hat{x}_{k+1} - \hat{x}_{k+1}^A + \hat{x}_{k+1}^A - x_{k+1})^2] \\ &= \mathbb{E}[(\hat{x}_{k+1} - \hat{x}_{k+1}^A)^2] + \mathbb{E}[(\hat{x}_{k+1}^A - x_{k+1})^2] \\ &\quad + 2\mathbb{E}[(\hat{x}_{k+1} - \hat{x}_{k+1}^A)(\hat{x}_{k+1}^A - x_{k+1})] \\ &= \mathbb{E}[(\tilde{e}_{k+1})^2] + P \\ &= \frac{\sigma_z^2}{c^2}(\bar{\delta}(\epsilon) - 1) + P \\ &= \bar{\delta}(\epsilon)P + \frac{(\bar{\delta}(\epsilon) - 1)\sigma_v^2}{c^2}. \end{aligned} \quad (24)$$

In (23), the fact  $\mathbb{E}[(\hat{x}_{k+1} - \hat{x}_{k+1}^A)(\hat{x}_{k+1}^A - x_{k+1})] = 0$  is due to the principle of orthogonality, i.e., all the random variables generated by  $\mathcal{I}_k$  is independent of the estimation error  $(\hat{x}_{k+1}^A - x_{k+1})$  of the MMSE estimate. Hence, the upper bound of  $\tilde{P}$  in (9) is achieved by this attack.

Now we show that the attack  $\tilde{u}_1^\infty$  is  $\epsilon$ -stealthy. From (21) and (22), we obtain  $\tilde{z}_k = z_k^A + c\zeta_{k-1}$ . Since  $\{z_k^A\}_{k=1}^\infty$  is an i.i.d. random sequence with  $z_k^A \sim \mathcal{N}(0, \sigma_z^2)$ , the random sequence  $\tilde{z}_1^\infty$  is i.i.d. Gaussian with  $\tilde{z}_k \sim \mathcal{N}(0, \bar{\delta}(\epsilon)\sigma_z^2)$ . For

every  $k > 0$ , we can calculate the KLD as

$$\begin{aligned} \frac{1}{k} \sum_{n=1}^k D(\tilde{y}_1^k \| y_1^k) &= \frac{1}{k} \sum_{n=1}^k D(\tilde{z}_1^k \| z_1^k) \\ &= \frac{1}{k} \sum_{n=1}^k -\frac{1}{2} \log(2\pi e \bar{\delta}(\epsilon) \sigma_z^2) + \frac{1}{2} \log(2\pi \sigma_z^2) + \frac{\bar{\delta}(\epsilon) \sigma_z^2}{2\sigma_z^2} \\ &= -\frac{1}{2} - \frac{1}{2} \log \bar{\delta}(\epsilon) + \frac{1}{2} \bar{\delta}(\epsilon) \\ &= \epsilon \end{aligned}$$

where the differential entropy of  $\tilde{z}_1^k$  is given by  $h(\tilde{z}_1^k) = \sum_{n=1}^k h(\tilde{z}_n) = \frac{k}{2} \log(2\pi e \bar{\delta}(\epsilon) \sigma_z^2)$  because  $\tilde{z}_1^\infty$  is an i.i.d. Gaussian sequence. In this case,  $\tilde{y}_1^\infty$  is ergodic. From Lemma 1, the attack  $\tilde{u}_1^\infty$  is  $\epsilon$ -stealthy. To conclude the proof, notice that the attack (20) can be generated by any information pattern satisfying (A1)–(A3). ■

**Remark 2: (Attacker information pattern)** As a counter-intuitive fact, Theorem 1 and Theorem 2 imply that knowledge of the system state does not increase the performance degradation induced by an attacker. In fact, the only critical piece of information for the attacker is the nominal control input  $u_1^\infty$ . It should be also noticed that knowledge of the nominal control input may not be necessary for different system and attack models. For instance, in the case the control input is transmitted via an additive channel, the attacker may achieve the upper bound (9) exploiting the linearity of the system, and without knowing the nominal control input. □

**Remark 3: (Properties of the optimal attack)** Recall that we make no assumption on the form of attacks. Yet, Theorem 2 implies that the random sequence  $\tilde{z}_1^\infty$  generated the optimal attack remains i.i.d. Gaussian with zero mean. This property follows from the fact that the inequalities (12), (13) and (14) hold with equalities in the case of optimal attacks.

## V. NUMERICAL RESULTS

We now present numerical results to illustrate the fundamental performance bounds derived in Section IV. The following results are stated based on the ratio  $\tilde{P}/P$ , which can be interpreted as the attacker gain. If the ratio  $\tilde{P}/P = 1$ , then the attacker can induce no degradation of the mean squared error. In Theorem 1 and Theorem 2 we characterize how an attacker must compromise between stealthiness and performance degradation at the system level. To illustrate such a trade-off, in Fig. 2 we report the ratio  $\tilde{P}/P$  as a function of the attack stealthiness  $\epsilon$ , for given system parameters.

In Fig. 3 we illustrate the relation between the attacker gain  $\tilde{P}/P$  and the quality of the measurements, as measured by  $c^2/\sigma_v^2$ . As expected, for a desired level of stealthiness, the attacker gain is smaller for larger values of  $c^2/\sigma_v^2$ .

Consider now the limiting situation of an unstable system with  $c^2/\sigma_v^2 \rightarrow 0^+$ . In this case the open loop unstable system is not detectable and thus  $P \rightarrow \infty$ . By taking the limit of (9) as  $c^2/\sigma_v^2 \rightarrow 0^+$  we obtain  $\tilde{P} \rightarrow \infty$ . In accordance with

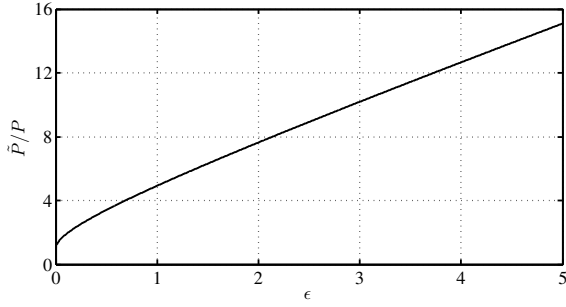


Fig. 2. This figure shows that attack stealthiness ( $\epsilon$ ) and performance degradation at the system level ( $\tilde{P}/P$ ) are competing objectives. The degradation  $\tilde{P}$  is induced by the optimal  $\epsilon$ -stealthy attack in (20). The system parameters are  $a = 2$ ,  $c = 1$ ,  $\sigma_w^2 = 0.5$ , and  $\sigma_v^2 = 0.1$ .

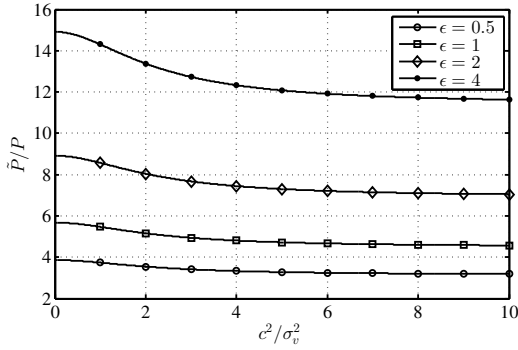


Fig. 3. This figure shows that, for a desired value of stealthiness, the larger the quality of measurements ( $c^2/\sigma_v^2$ ) the smaller the attacker gain ( $\tilde{P}/P$ ). The system parameters are  $a = 2$  and  $\sigma_w^2 = 0.5$ , and the degradation  $\tilde{P}$  is induced by the optimal  $\epsilon$ -stealthy attack in (20).

these results, Fig. 3 shows that  $\tilde{P}/P$  remains bounded as  $c^2/\sigma_v^2 \rightarrow 0^+$ .

Similarly, we consider the limiting situation of a stable system with  $c^2/\sigma_v^2 \rightarrow 0^+$ . The attacker gain  $\tilde{P}/P$  as a function of  $c^2/\sigma_v^2$  is reported in Fig. 4. It can be observed that  $\tilde{P}/P$  grows unbounded as  $c^2/\sigma_v^2 \rightarrow 0^+$ . In fact, since the system is stable, the mean squared error of the Kalman filter  $P$  is bounded for all  $c^2/\sigma_v^2 \geq 0$ . On the other hand, by taking the limit of (9) we observe that  $\tilde{P}$  goes to infinity as  $c^2/\sigma_v^2 \rightarrow 0^+$ .

## VI. CONCLUSION

This work characterizes fundamental limitations and performance bounds for the security of stochastic control systems. The scenario is considered where the attacker knows the system parameters and noise statistics, and is able to hijack and replace the nominal control input. We propose a notion of  $\epsilon$ -stealthiness to quantify the difficulty to detect an attack from measurements, and we characterize the maximal degradation of the control performance induced by an  $\epsilon$ -stealthy attack. Our study reveals that an  $\epsilon$ -stealthy attacker only need to know the nominal control input to cause the largest performance degradation in Kalman filtering.

## REFERENCES

[1] A. Teixeira, D. Pérez, H. Sandberg, and K. H. Johansson, "Attack models and scenarios for networked control systems," in *Proc. of the*

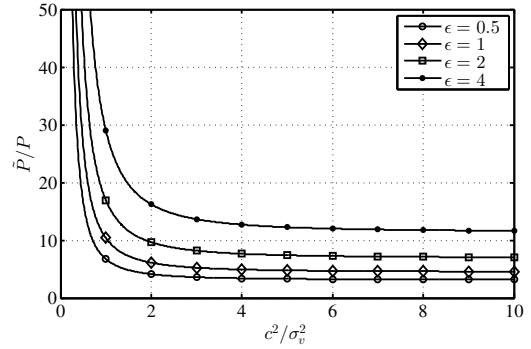


Fig. 4. This figure shows the tradeoff between performance degradation at the system level ( $\tilde{P}/P$ ) and the quality of measurements ( $c^2/\sigma_v^2$ ) for a stable system. The system parameters are  $a = 0.5$  and  $\sigma_w^2 = 0.5$ , and the degradation  $\tilde{P}$  is induced by the optimal  $\epsilon$ -stealthy attack in (20). Notice that, contrarily to the case of unstable system in Fig. 3, the attacker gain grows unbounded as  $c^2/\sigma_v^2$  approaches zero.

*1st international conference on High Confidence Networked Systems*. ACM, 2012, pp. 55–64.

[2] H. S. Foroush and S. Martínez, "On multi-input controllable linear systems under unknown periodic dos jamming attacks," in *SIAM Conf. on Control and its Applications*. SIAM, 2013, pp. 222–229.

[3] R. S. Smith, "A decoupled feedback structure for covertly appropriating control systems," *Network*, vol. 6, p. 6, 2011.

[4] Y. Mo and B. Sinopoli, "Secure control against replay attacks," in *47th Annual Allerton Conference*. IEEE, 2009, pp. 911–918.

[5] O. Kosut, L. Jia, R. J. Thomas, and L. Tong, "Malicious data attacks on the smart grid," *Smart Grid, IEEE Trans. on*, vol. 2, no. 4, pp. 645–658, 2011.

[6] Y. Liu, P. Ning, and M. K. Reiter, "False data injection attacks against state estimation in electric power grids," *ACM Trans. on Information and System Security*, vol. 14, no. 1, p. 13, 2011.

[7] C. Kwon, W. Liu, and I. Hwang, "Security analysis for cyber-physical systems against stealthy deception attacks," in *American Control Conference (ACC), 2013*. IEEE, 2013, pp. 3344–3349.

[8] Y. Mo, R. Chabukswar, and B. Sinopoli, "Detecting integrity attacks on scada systems," *IEEE Transactions on Control Systems Technology*, vol. 22, no. 4, pp. 1396–1407, 2014.

[9] G. Basile and G. Marro, *Controlled and Conditioned Invariants in Linear System Theory*. Prentice Hall, 1991.

[10] F. Pasqualetti, F. Dörfler, and F. Bullo, "Attack detection and identification in cyber-physical systems," *IEEE Trans. Autom. Control*, vol. 58, no. 11, 2013.

[11] H. Fawzi, P. Tabuada, and S. Diggavi, "Secure estimation and control for cyber-physical systems under adversarial attacks," *IEEE Trans. on Automatic Control*, vol. 59, no. 6, pp. 1454–1467, 2014.

[12] S. Cui, Z. Han, S. Kar, T. T. Kim, H. V. Poor, and A. Tajer, "Coordinated data-injection attack and detection in the smart grid: A detailed look at enriching detection solutions," *Signal Processing Magazine, IEEE*, vol. 29, no. 5, pp. 106–115, 2012.

[13] C.-Z. Bai and V. Gupta, "On Kalman filtering in the presence of a compromised sensor: Fundamental performance bounds," in *American Control Conference (ACC)*, Portland, OR, June 2014, pp. 3029–3034.

[14] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Wiley, 2006.

[15] R. E. Kalman, "A new approach to linear filtering and prediction problems," *J. of Basic Engineering*, vol. 82, no. 1, pp. 35–45, 1960.

[16] T. Kailath, A. H. Sayed, and B. Hassibi, *Linear Estimation*. Prentice Hall, 2000.

[17] H. V. Poor, *An introduction to signal detection and estimation*, 2nd ed. New York: Springer-Verlag, 1998.

[18] S. Kullback, *Information theory and statistics*. Courier Dover, 1997.

[19] Y. Polyanskiy and Y. Wu, *Lecture notes on Information Theory*. MIT (6.441), UIUC (ECE 563), 2012–2013.

[20] C.-Z. Bai, F. Pasqualetti, and V. Gupta, "Notes on security in stochastic control systems: Fundamental limitations and performance bounds (ACC 2015)," <http://www3.nd.edu/~vgupta2/research/publications/ACC2015Note.pdf>.