

# Kalman Filter Based Secure State Estimation and Individual Attacked Sensor Detection in Cyber-Physical Systems[1]

Presented by: M. H. Basiri and J. G. Thistle and J. W. Simpson-Porco and S. Fischmeister  
at the 2019 American Control Conference (ACC)

Jonas Wagner

MECH 6325 Research Seminar  
Fall 2020

# Outline

- 1 Overview
- 2 System and Attack Modeling
- 3 Attack Detection Preliminaries
  - Kalman Filtering
  - Conventional  $\chi^2$ -Detector
- 4 Proposed Attack Detection Algorithms
  - RWD Attack Detector
  - NRD Attack Detector
- 5 Simulation Result
  - The Unstealthy Case
  - The Stealthy Case
  - The Very Unstealthy Case
  - Detection Rate Analysis

# Outline

- 1 Overview**
- 2 System and Attack Modeling**
- 3 Attack Detection Preliminaries**
  - Kalman Filtering
  - Conventional  $\chi^2$ -Detector
- 4 Proposed Attack Detection Algorithms**
  - RWD Attack Detector
  - NRD Attack Detector
- 5 Simulation Result**
  - The Unstealthy Case
  - The Stealthy Case
  - The Very Unstealthy Case
  - Detection Rate Analysis

# Paper Information

**Title:** Kalman Filter Based Secure State Estimation and Individual Attacked Sensor Detection in Cyber-Physical Systems

**Authors:** M. H. Basiri and J. G. Thistle and J. W. Simpson-Porco and S. Fischmeister

**Conference:** 2019 American Control Conference (ACC)

**Main Contributions:** Development and simulation of two attack detection and secure state estimation algorithms - Rolling Window Detector (RWD) and Novel Residual Detector (NRD)

# Security of Cyber-Physical Systems

- Cyber-Physical Systems (CPS) allow for large scale wide spread control systems that can take advantages of wireless data communication
- CPSs are unfortunately vulnerable to malicious attacks
- Attacks on a subset of sensors may attempt to disrupt the performance of such system
- Detection of such attacks is important and it is difficult to distinguish between noise and attacks
- Many detection methods actually use Kalmen filter estimation and perform hypothetical testing base the residual vector to detect changes above the standard noise threshold

# Outline

- 1 Overview
- 2 System and Attack Modeling
- 3 Attack Detection Preliminaries
  - Kalman Filtering
  - Conventional  $\chi^2$ -Detector
- 4 Proposed Attack Detection Algorithms
  - RWD Attack Detector
  - NRD Attack Detector
- 5 Simulation Result
  - The Unstealthy Case
  - The Stealthy Case
  - The Very Unstealthy Case
  - Detection Rate Analysis

# System and Attack Definitions

## Original System

$$\mathcal{P} : \begin{cases} x_{k+1} = Ax_k + Bu_k + \nu_k, \\ y_k = Cx_k + w_k \end{cases} \quad (1)$$

**States, Inputs, and Outputs:**

$x \in \mathbb{R}^n$ ,  $u \in \mathbb{R}^m$ , and  $y \in \mathbb{R}^q$

**Noise:**  $\nu \sim \mathcal{N}(0, Q)$  and  $w \sim \mathcal{N}(0, R)$

**Initial State:**  $x_0 \sim \mathcal{N}(0, \Sigma)$

## Attacked System

$$\mathcal{P}_a : \begin{cases} x_{k+1}^a = Ax_k^a + Bu_k + \nu_k, \\ y_k^a = Cx_k^a + w_k + Da_k \end{cases} \quad (2)$$

**States, Inputs, and Outputs:**  $x_k^a$ ,  $u_k$ , and  $y_k^a$

**Noise:**  $\nu \sim \mathcal{N}(0, Q)$  and  $w \sim \mathcal{N}(0, R)$

**Attacks:**  $a_k$  is manipulated

$D$  s.t.  $D_{jj} = 1$  if  $j^{th}$  sensor is under attack, otherwise  $D_{ij} = 0 \forall i, j$

# Outline

- 1 Overview
- 2 System and Attack Modeling
- 3 Attack Detection Preliminaries
  - Kalman Filtering
  - Conventional  $\chi^2$ -Detector
- 4 Proposed Attack Detection Algorithms
  - RWD Attack Detector
  - NRD Attack Detector
- 5 Simulation Result
  - The Unstealthy Case
  - The Stealthy Case
  - The Very Unstealthy Case
  - Detection Rate Analysis

# Kalman Filter Definition

## Measurement Update

$$K_k = P_{k|k-1} C^T (C P_{k|k-1} C^T + R)^{-1} \quad (3)$$

$$P_{k|k} = P_{k|k-1} - K_k C P_{k|k-1} \quad (4)$$

$$\hat{x}_{k|k} = \hat{x}_{k|k-1} + K_k (y_k - C \hat{x}_{k|k-1}) \quad (5)$$

## Equivalent Equations:

$$K_k = P_k^- H^T (H P_k^- H^T + R)^{-1}$$

$$P_k^+ = (I - K_k H_k) P_k^-$$

$$\hat{x}_k^+ = \hat{x}_k^- + K_k (y_k - H \hat{x}_k^-)$$

## Time Update

$$\hat{x}_{k+1|k} = A \hat{x}_{k|k} + B u_k \quad (6)$$

$$P_{k+1|k} = A P_{k|k} A^T + Q \quad (7)$$

$$\hat{x}_k^- = F \hat{x}_{k-1}^+ + G u_{k-1}$$

$$P_k^- = F P_{k-1}^+ F^T + Q$$

**Steady-State:**  $P = \triangle \lim_{k \rightarrow \infty} P_{k|k-1}$      $K = \triangle \lim_{k \rightarrow \infty} K_k = P C^T (C P C^T + R)^{-1}$

# Conventional $\chi^2$ -Detector

Widely used in fault-detection and also applicable to CPS security.

**Residual Vector:**

$$r_k = y_k - \hat{y}_k = y_k - C\hat{x}_k \quad (8)$$

**Residual Covariance:**

$$\Sigma_{r,k} = Ck|kC^T + R \quad (9)$$

**Residual Power:**

$$g_k = r_k^T \Sigma_{r,k}^{-1} r_k \quad (10)$$

**Detector Statistical Test:**

A threshold determined using the desired confidence and degrees of freedom from the *chi*<sup>2</sup>-distribution. An alarm is then triggered if:

$$g_k > \text{threshold} \quad (11)$$

# Outline

- 1 Overview
- 2 System and Attack Modeling
- 3 Attack Detection Preliminaries
  - Kalman Filtering
  - Conventional  $\chi^2$ -Detector
- 4 Proposed Attack Detection Algorithms
  - RWD Attack Detector
  - NRD Attack Detector
- 5 Simulation Result
  - The Unstealthy Case
  - The Stealthy Case
  - The Very Unstealthy Case
  - Detection Rate Analysis

# Rolling Window Detector (RWD)

This detector works by comparing the measured residual cumulative sum  $\bar{P}_{k|k}$  over a set period of time  $T$  with the error covariance predicted by the Kalman filter,  $P_{k|k}$ .

## Test Definitions:

$$\hat{\Sigma}_k = \frac{1}{T} \sum_{k=k_0}^{k_0+T-1} (y_k - \hat{y}_k)(y_k - \hat{y}_k)^T \quad (12)$$

$$\Sigma_{r,k} = CP_{k|k}C^T + R \quad (13)$$

**Testing Statistic:** ( $H$  is the threshold matrix)

$$S(T, k_0) = \hat{\Sigma}_k - \Sigma_{r,k} > H \quad (14)$$

---

### 1: Initialize:

$x_0 \sim \mathcal{N}(0, \Sigma)$ ,  $\hat{x}_{0|-1} = \mathbf{0}$ ,  $P_{0|-1} = \Sigma$ ,  $Q$ ,  $R$ ,  $T$ ,  $H$ .

2: **while**  $k \leq k_{\text{final}}$  **do**

3:    Make a backup of  $\hat{x}_{k|k-1}$  in  $\hat{x}_{k,\text{backup}}$ .

4:    **for** Each sensor  $j$  **do**

5:     Apply standard Kalman filter and estimate  $\hat{x}_{k|k}$

6:     **calculate:**  $\hat{\Sigma}_k = \frac{1}{T} \sum_{k_0}^{k_0+T-1} (y_k - \hat{y}_k)(y_k - \hat{y}_k)^T$

7:     **calculate:**  $\Sigma_{r,k} = CP_{k|k}C^T + R$

8:     **if**  $(\hat{\Sigma}_k - \Sigma_{r,k})_{jj} > H_{jj}$  **then**

9:       Modify  $j^{\text{th}}$  component of vector  $y_k$  by replacing it with the  $j^{\text{th}}$  component of vector  $\hat{y}_{k-1}$ .

10:      Re-estimate  $\hat{x}_{k|k}$  via (6) using  $\hat{x}_{k,\text{backup}}$  and the new  $y_k$ .

11:      Re-estimate  $\hat{y}_k$  with the new  $\hat{x}_{k|k}$ .

12:     **end if**

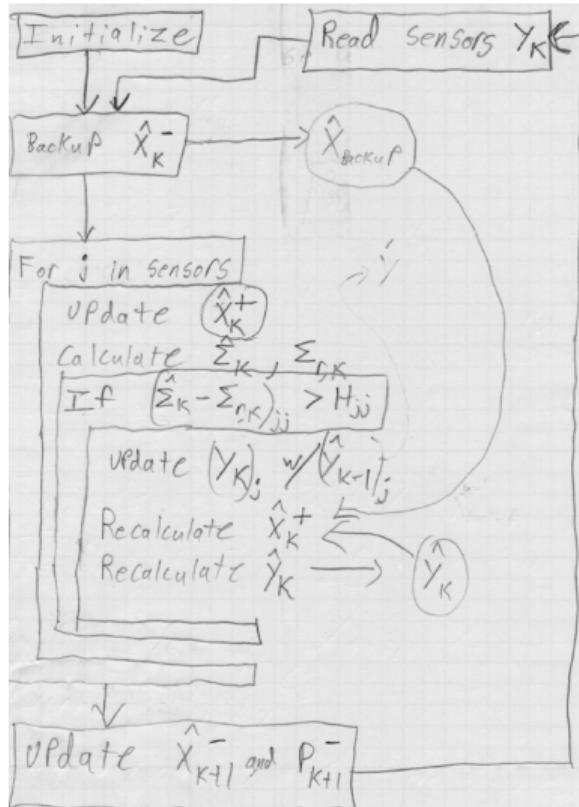
13:    **end for**

14:    Predict  $\hat{x}_{k+1|k}$  and  $P_{k+1|k}$  via (7) and (8).

15: **end while**

---

# RWD Procedure



**1: Initialize:**

$x_0 \sim \mathcal{N}(0, \Sigma)$ ,  $\hat{x}_{0|-1} = \mathbf{0}$ ,  $P_{0|-1} = \Sigma$ ,  $Q$ ,  $R$ ,  $T$ ,  $H$ .

**2: while**  $k \leq k_{\text{final}}$  **do**

**3:** Make a backup of  $\hat{x}_{k|k-1}$  in  $\hat{x}_{k,\text{backup}}$ .

**4:** **for** Each sensor  $j$  **do**

**5:** Apply standard Kalman filter and estimate  $\hat{x}_{k|k}$

**6:** **calculate:**  $\hat{\Sigma}_k = \frac{1}{T} \sum_{k_0}^{k_0+T-1} (y_k - \hat{y}_k)(y_k - \hat{y}_k)^T$

**7:** **calculate:**  $\Sigma_{r,k} = CP_{k|k}C^T + R$

**8:** **if**  $(\hat{\Sigma}_k - \Sigma_{r,k})_{jj} > H_{jj}$  **then**

**9:** Modify  $j^{\text{th}}$  component of vector  $y_k$  by replacing it with the  $j^{\text{th}}$  component of vector  $\hat{y}_{k-1}$ .

**10:** Re-estimate  $\hat{x}_{k|k}$  via (6) using  $\hat{x}_{k,\text{backup}}$  and the new  $y_k$ .

**11:** Re-estimate  $\hat{y}_k$  with the new  $\hat{x}_{k|k}$ .

**12:** **end if**

**13:** **end for**

**14:** Predict  $\hat{x}_{k+1|k}$  and  $P_{k+1|k}$  via (7) and (8).

**15:** **end while**

# Novel Residual Detector (NRD)

This method extends the  $\chi^2$ -detector method to individual sensors using a similar modified Kalman filter to RWD.

**Test Definitions:** (First is identical to  $\chi^2$ -test)

$$g_k = r_k^T \Sigma_{r,k}^{-1} r_k > \text{threshold}_1 \quad (15)$$

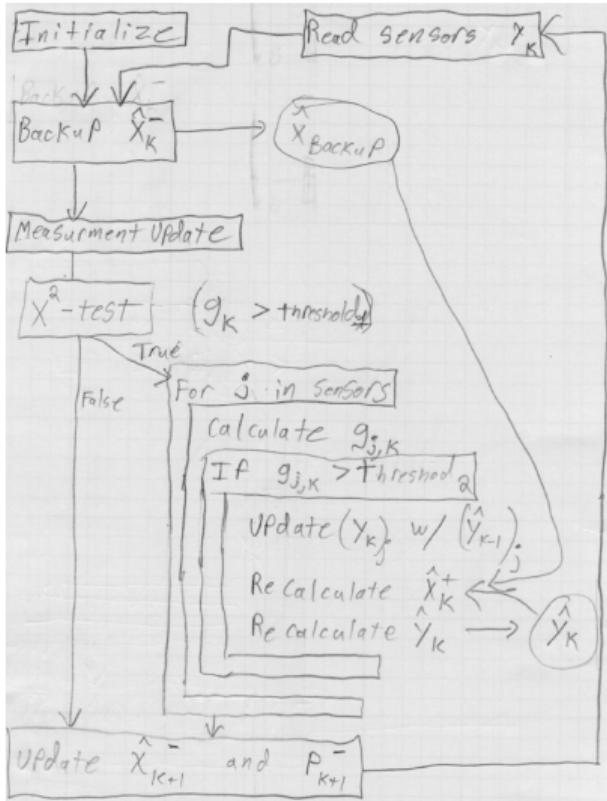
$$g_{j,k} = \frac{r_k^2}{(\Sigma_{r,k})_{jj}} > \text{threshold}_2 \quad (16)$$

---

```

1: Initialize:  $x_0 \sim \mathcal{N}(0, \Sigma)$ ,  $\hat{x}_{0|-1} = \mathbf{0}$ ,  $P_{0|-1} = \Sigma$ ,  $Q$ ,  $R$ .
2: while  $k \leq k_{\text{final}}$  do
3:   Make a backup of  $\hat{x}_{k|k-1}$  in  $\hat{x}_{k,\text{backup}}$ .
4:   Apply standard Kalman filter and estimate  $\hat{x}_{k|k}$ 
5:   calculate:  $g_k = r_k^T \Sigma_{r,k}^{-1} r_k$ 
6:   if  $g_k = r_k^T \Sigma_{r,k}^{-1} r_k > \text{threshold}_1$  then
7:     for Each sensor  $j$  do
8:       calculate:  $g_{j,k} = r_k^2 / (\Sigma_{r,k})_{jj}$ 
9:       if  $g_{j,k} = r_k^2 / (\Sigma_{r,k})_{jj} > \text{threshold}_2$  then
10:        Modify  $j^{\text{th}}$  component of vector  $y_k$  by
          replacing it with the  $j^{\text{th}}$  component of vector  $\hat{y}_{k-1}$ .
11:        Re-estimate  $\hat{x}_{k|k}$  via (6) using  $\hat{x}_{k,\text{backup}}$ 
          and the new  $y_k$ .
12:        Re-estimate  $\hat{y}_k$  with the new  $\hat{x}_{k|k}$ .
13:      end if
14:    end for
15:  end if
16:  Predict  $\hat{x}_{k+1|k}$  and  $P_{k+1|k}$  via (7) and (8).
17: end while
```

# NRD Procedure


**1: Initialize:**
 $x_0 \sim \mathcal{N}(0, \Sigma), \hat{x}_{0|-1} = \mathbf{0}, P_{0|-1} = \Sigma, Q, R.$ 
**2: while**  $k \leq k_{\text{final}}$  **do**

 3: Make a backup of  $\hat{x}_{k|k-1}$  in  $\hat{x}_{k,\text{backup}}$ .

 4: Apply standard Kalman filter and estimate  $\hat{x}_{k|k}$ 

 5: **calculate:**  $g_k = r_k^T \Sigma_{r,k}^{-1} r_k$ 

 6: **if**  $g_k = r_k^T \Sigma_{r,k}^{-1} r_k > \text{threshold}_1$  **then**

 7:   **for** Each sensor  $j$  **do**

 8:     **calculate:**  $g_{j,k} = r_k^2 / (\Sigma_{r,k})_{jj}$ 

 9:     **if**  $g_{j,k} = r_k^2 / (\Sigma_{r,k})_{jj} > \text{threshold}_2$  **then**

 10:       Modify  $j^{\text{th}}$  component of vector  $y_k$  by  
replacing it with the  $j^{\text{th}}$  component of vector  $\hat{y}_{k-1}$ .

 11:       Re-estimate  $\hat{x}_{k|k}$  via (6) using  $\hat{x}_{k,\text{backup}}$   
and the new  $y_k$ .

 12:       Re-estimate  $\hat{y}_k$  with the new  $\hat{x}_{k|k}$ .

 13:     **end if**

 14:   **end for**

 15: **end if**

 16: Predict  $\hat{x}_{k+1|k}$  and  $P_{k+1|k}$  via (7) and (8).

 17: **end while**

# Outline

- 1 Overview
- 2 System and Attack Modeling
- 3 Attack Detection Preliminaries
  - Kalman Filtering
  - Conventional  $\chi^2$ -Detector
- 4 Proposed Attack Detection Algorithms
  - RWD Attack Detector
  - NRD Attack Detector
- 5 Simulation Result
  - The Unstealthy Case
  - The Stealthy Case
  - The Very Unstealthy Case
  - Detection Rate Analysis

## IEEE 14-bus Power Grid System [2]



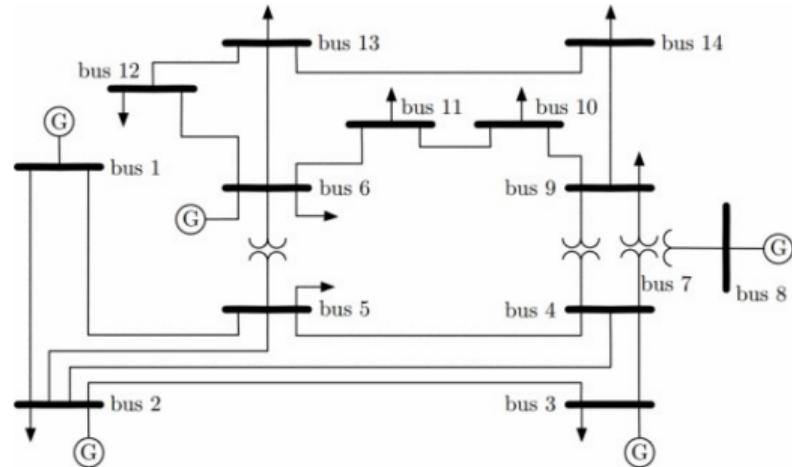
Effectiveness of RWD and NRD is tested on the IEEE 14-bus power grid system

- 5 generators and 14 buses
  - 14 power injection and 20 power flow sensors

**System Parameters:**  $n = 10$  and  $p = 35$

## System States:

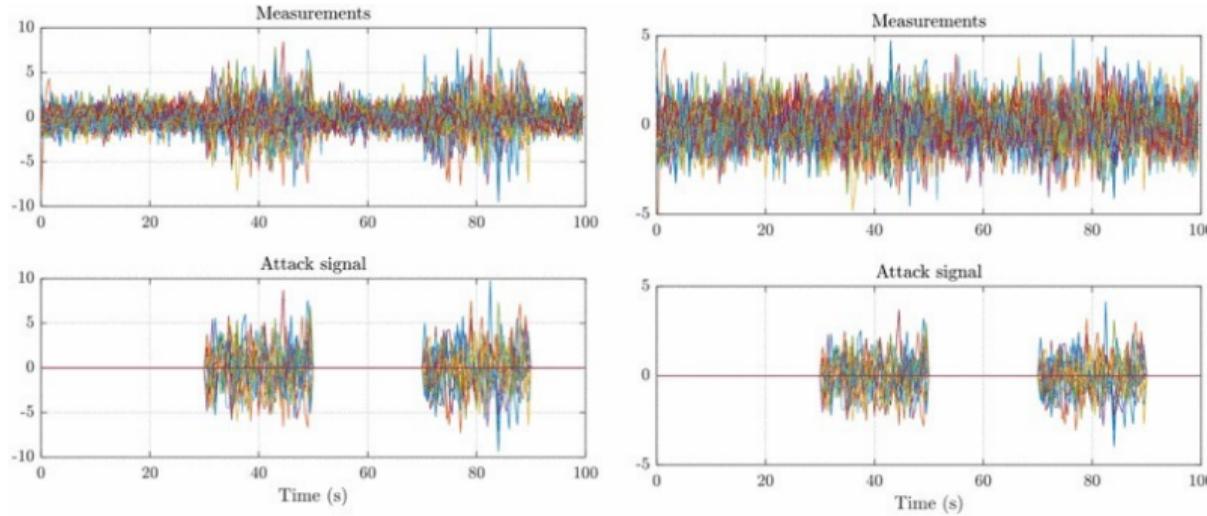
$\delta_i$  (rotor angle),  $\omega_i$  (frequencies)  $\forall i = \{1, \dots, \frac{n}{2}\}$



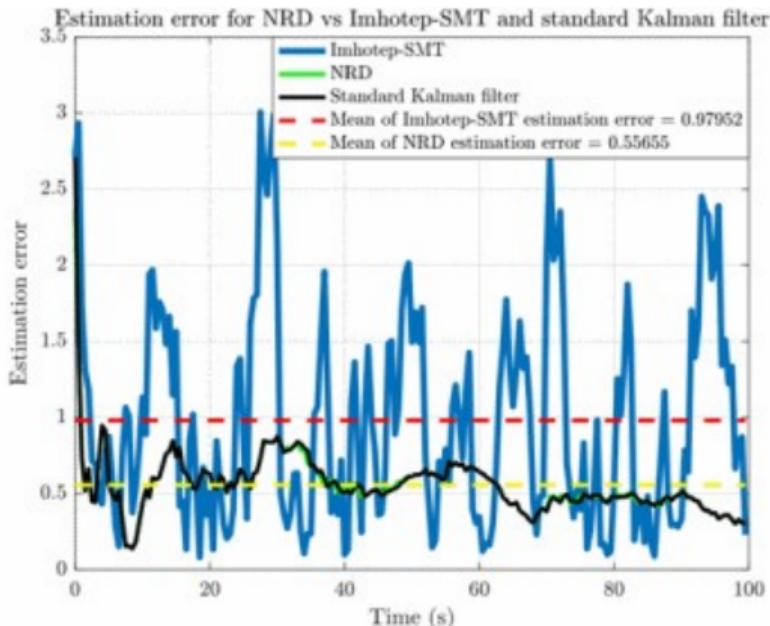
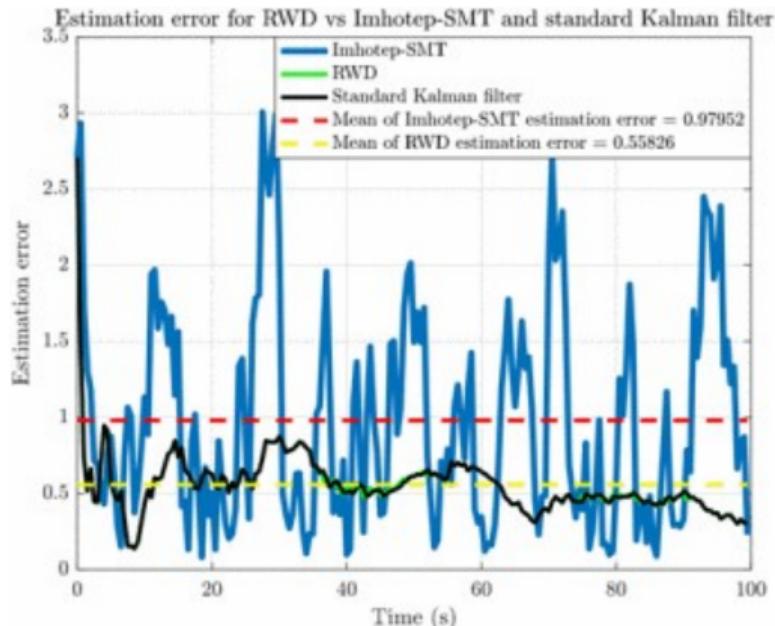
# Simulation Results Overview

The simulation tested the two methods for 3 cases:

- The Unstealthy Case - Attempts to disrupt system without remaining undetected
- The Stealthy Case - Injects false data upon the same order as the noise
- The Very Unstealthy Case - An unstealthy attack much larger than the noise

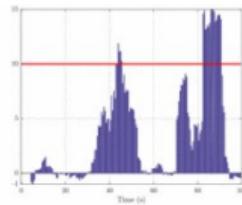


# The Unsteady Case

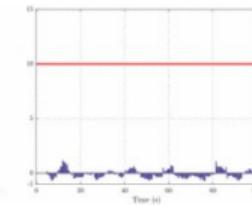


Estimation error compared with the Imhotep-SMT tool [3] and a standard Kalman Filter.

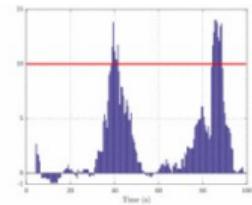
# The Unsteady Case - Individual Sensors



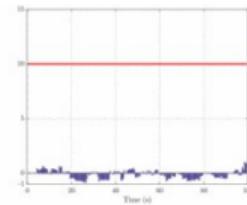
(a) Compromised sensor 8 with RWD



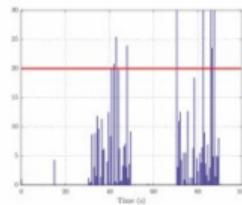
(b) Uncompromised sensor 9 with RWD



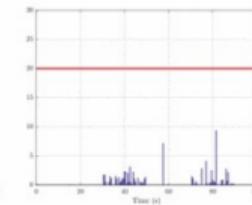
(c) Compromised sensor 15 with RWD



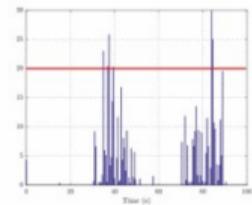
(d) Uncompromised sensor 34 with RWD



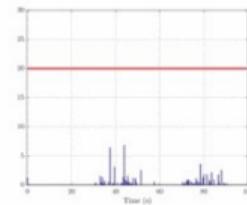
(e) Compromised sensor 8 with NRD



(f) Uncompromised sensor 9 with NRD



(g) Compromised sensor 15 with NRD

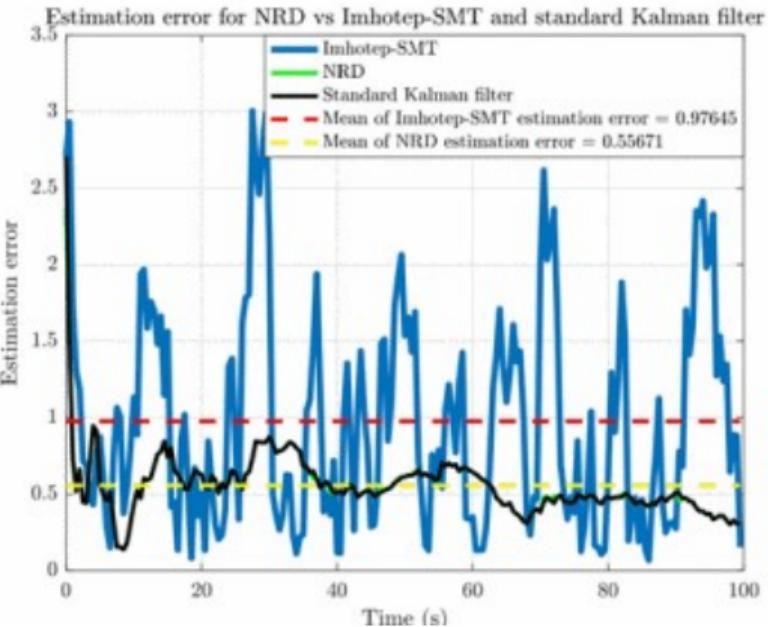
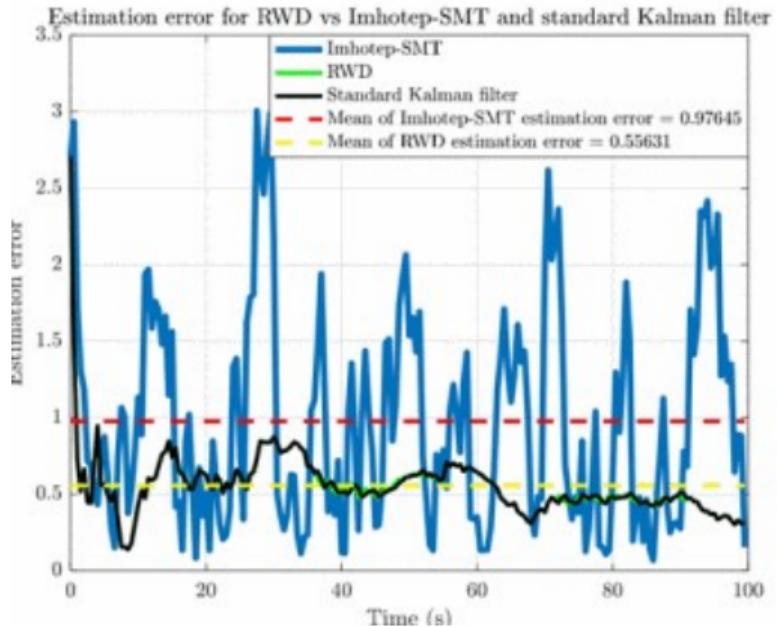


(h) Uncompromised sensor 34 with NRD

Residuals calculated for individual sensors for RWD and NRD (unsteady attacks).

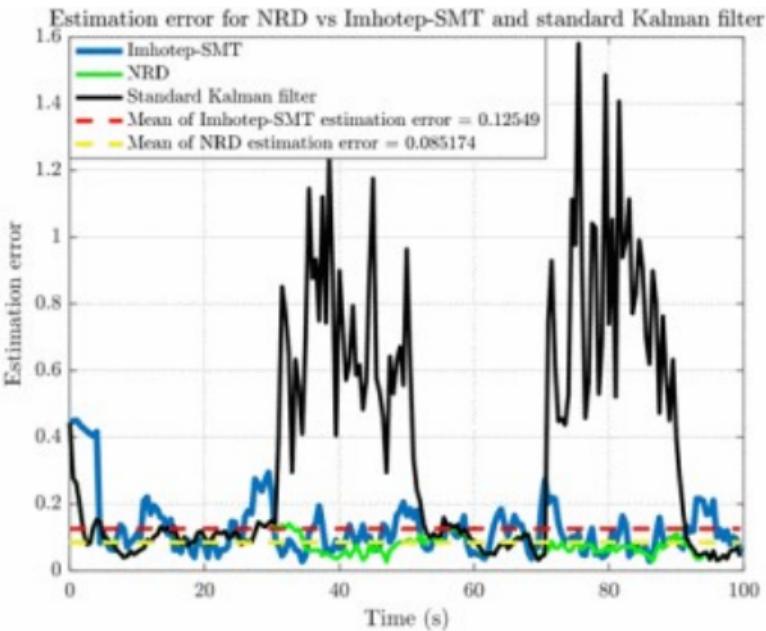
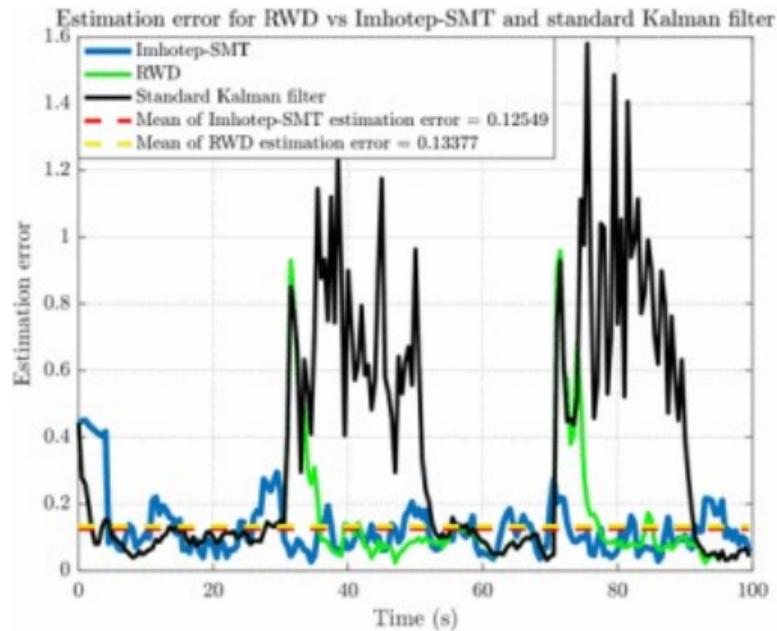


# The Stealthy Case



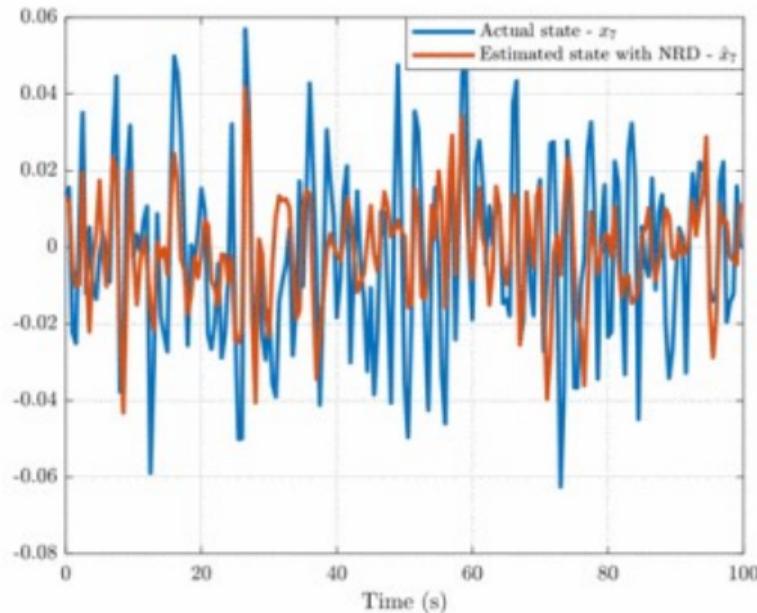
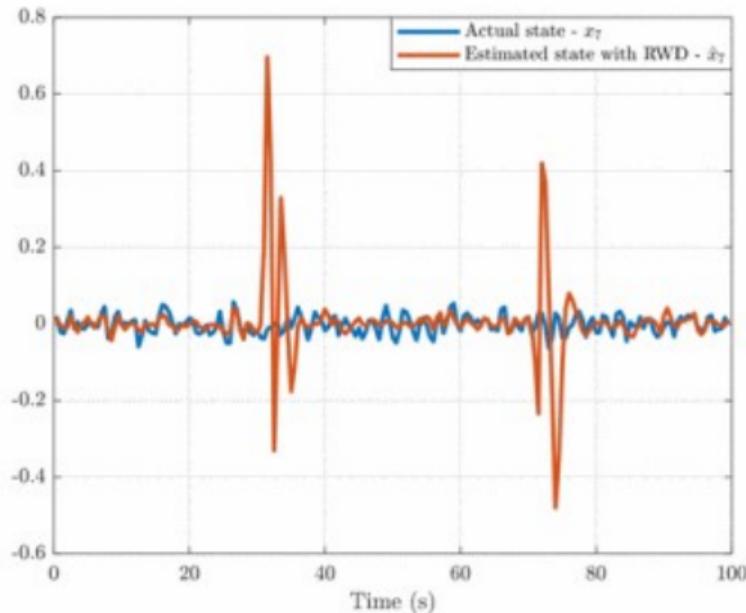
Estimation error compared with the Imhotep-SMT tool [3] and a standard Kalman Filter.

# The Very Unsteady Case



Estimation error compared with the Imhotep-SMT tool [3] and a standard Kalman Filter.

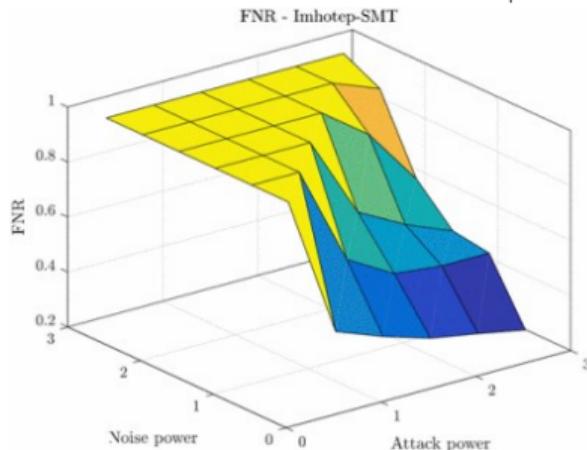
# The Very Unstealthy Case - Individual Sensor



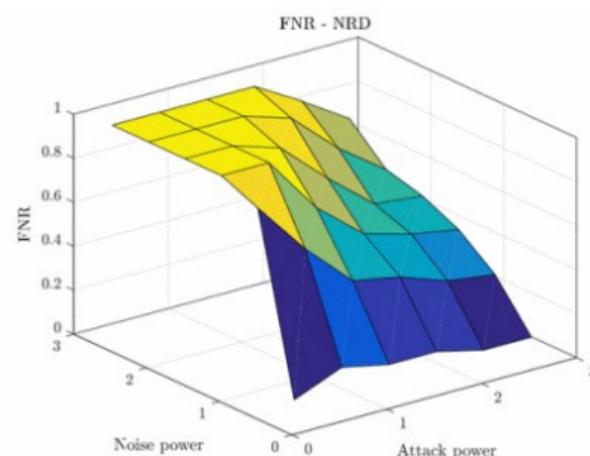
Estimation error of state  $x_7$  for RWD and NRD (very unstealthy attacks).

# Detection Rate Analysis

**False Positive Rate (FPR):**  $FNR = \frac{FN}{FP+TP}$



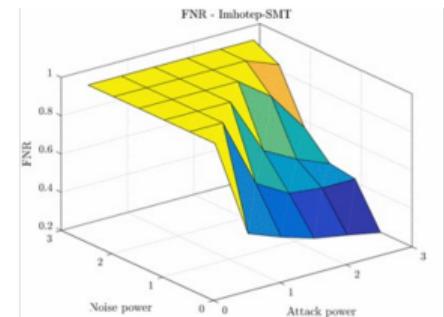
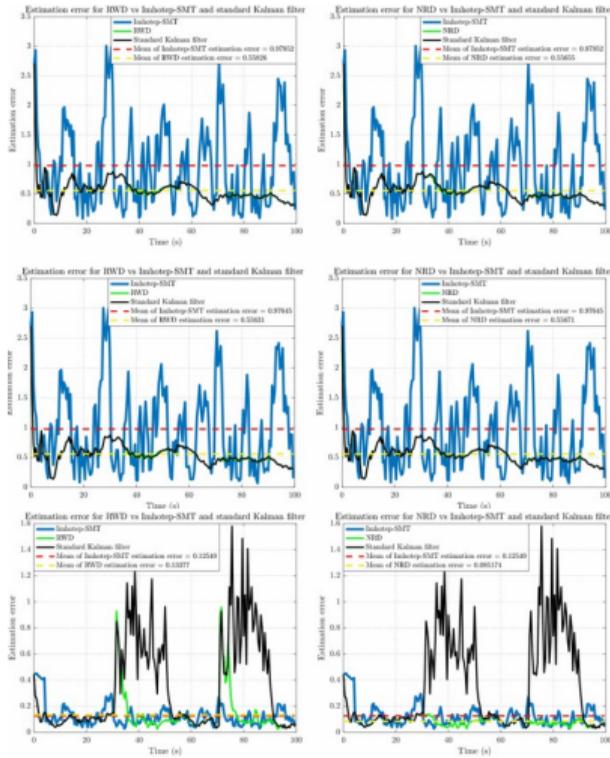
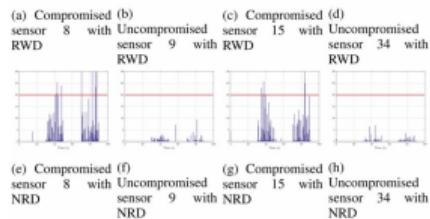
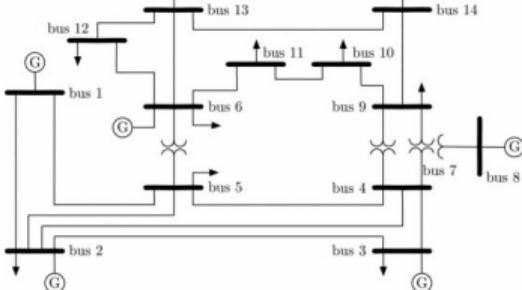
(a) FNR for Imhotep-SMT



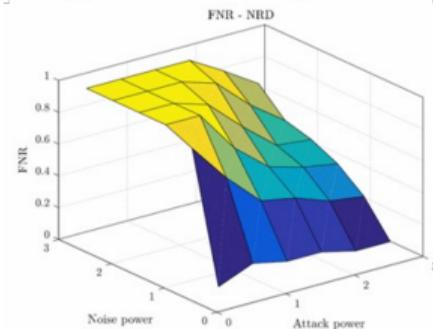
(b) FNR for NRD

FNR for Imhotep-SMT vs NRD. NRD clearly performs better than Imhotep-SMT when the attack and noise power decrease (aka stealthy case).

# Questions?



(a) FNR for Imhotep-SMT



(b) FNR for NRD

-  M. H. Basiri, J. G. Thistle, J. W. Simpson-Porco, and S. Fischmeister.  
Kalman filter based secure state estimation and individual attacked sensor detection in cyber-physical systems.  
In *2019 American Control Conference (ACC)*, pages 3841–3848, 2019.
-  F. Pasqualetti, F. Dörfler, and F. Bullo.  
Cyber-physical attacks in power networks: Models, fundamental limitations and monitor design.  
In *2011 50th IEEE Conference on Decision and Control and European Control Conference*, pages 2195–2201, 2011.

-  Y. Shoukry, P. Nuzzo, A. Puggelli, A. L. Sangiovanni-Vincentelli, S. A. Seshia, and P. Tabuada.  
Secure state estimation for cyber-physical systems under sensor attacks: A satisfiability modulo theory approach.  
*IEEE Transactions on Automatic Control*, 62(10):4917–4932, 2017.