



KTH ROYAL INSTITUTE  
OF TECHNOLOGY

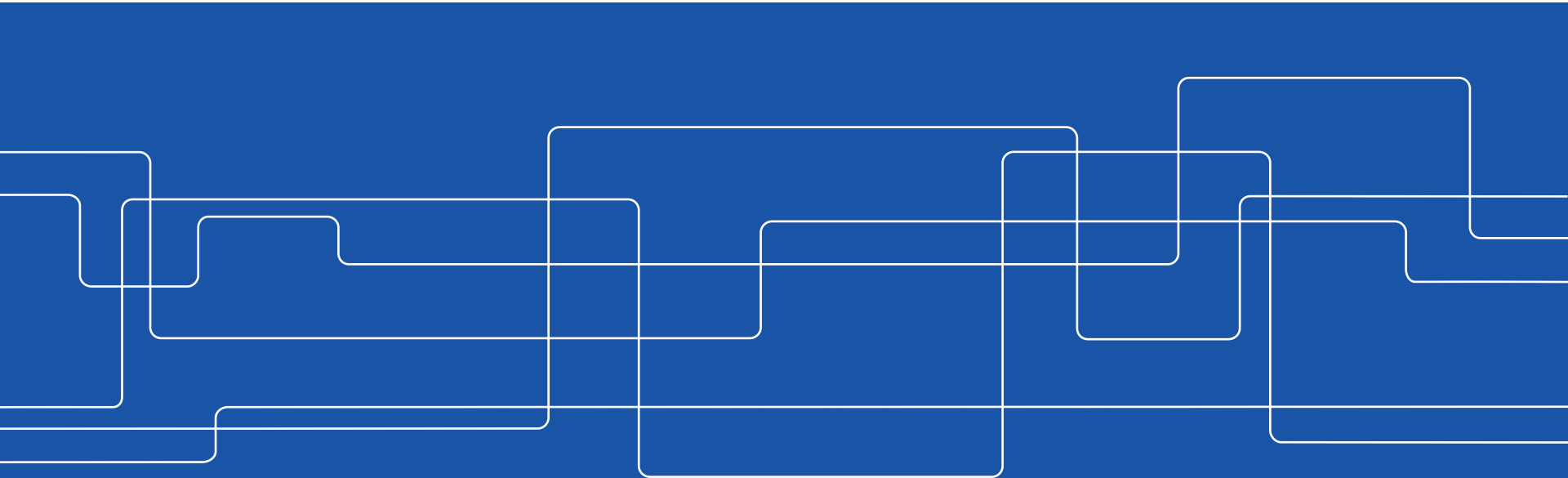
# Security of Cyber-Physical Systems

**Henrik Sandberg**

[hsan@kth.se](mailto:hsan@kth.se)

Department of Automatic Control, KTH, Stockholm, Sweden

7th oCPS PhD School on Cyber-Physical Systems, Lucca, Italy





# Outline

- Background and motivation
- CPS attack models
- Risk management
- Attack detectability and security metrics
- Attack identification and secure state estimation

# ICS-CERT MONITOR



January – April 2014



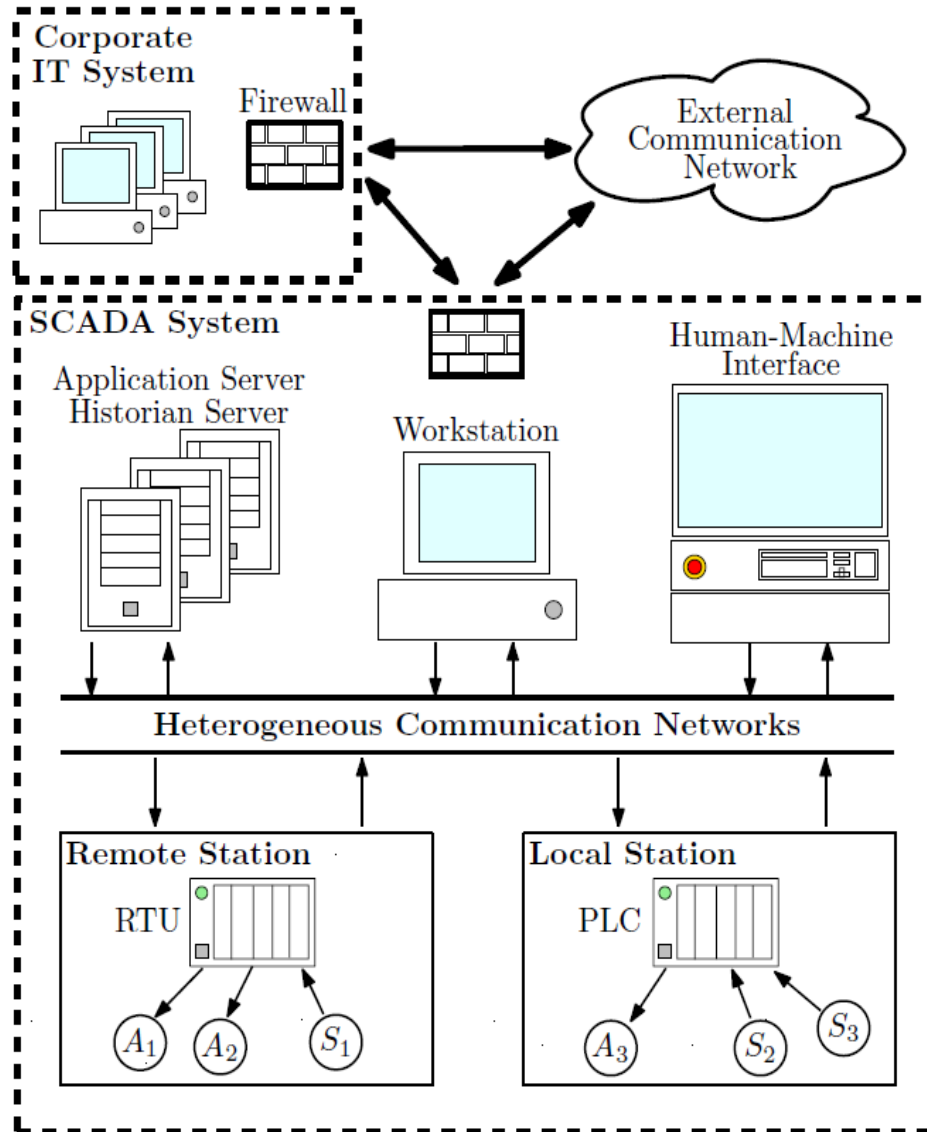
## INCIDENT RESPONSE ACTIVITY

### INTERNET ACCESSIBLE CONTROL SYSTEMS AT RISK

Is your control system accessible directly from the Internet? Do you use remote access features to log into your control system network? Are you unsure of the security measures that protect your remote access services? If your answer was yes to any or all these questions, you are at increased risk of cyber attacks including scanning, probes, brute force attempts and unauthorized access to your control environment.

ICS-CERT = Industrial Control Systems Cyber Emergency Response Team  
(<https://ics-cert.us-cert.gov/>)  
Part of US Department of Homeland Security

# Example 1: Industrial Control System (ICS) Infrastructure



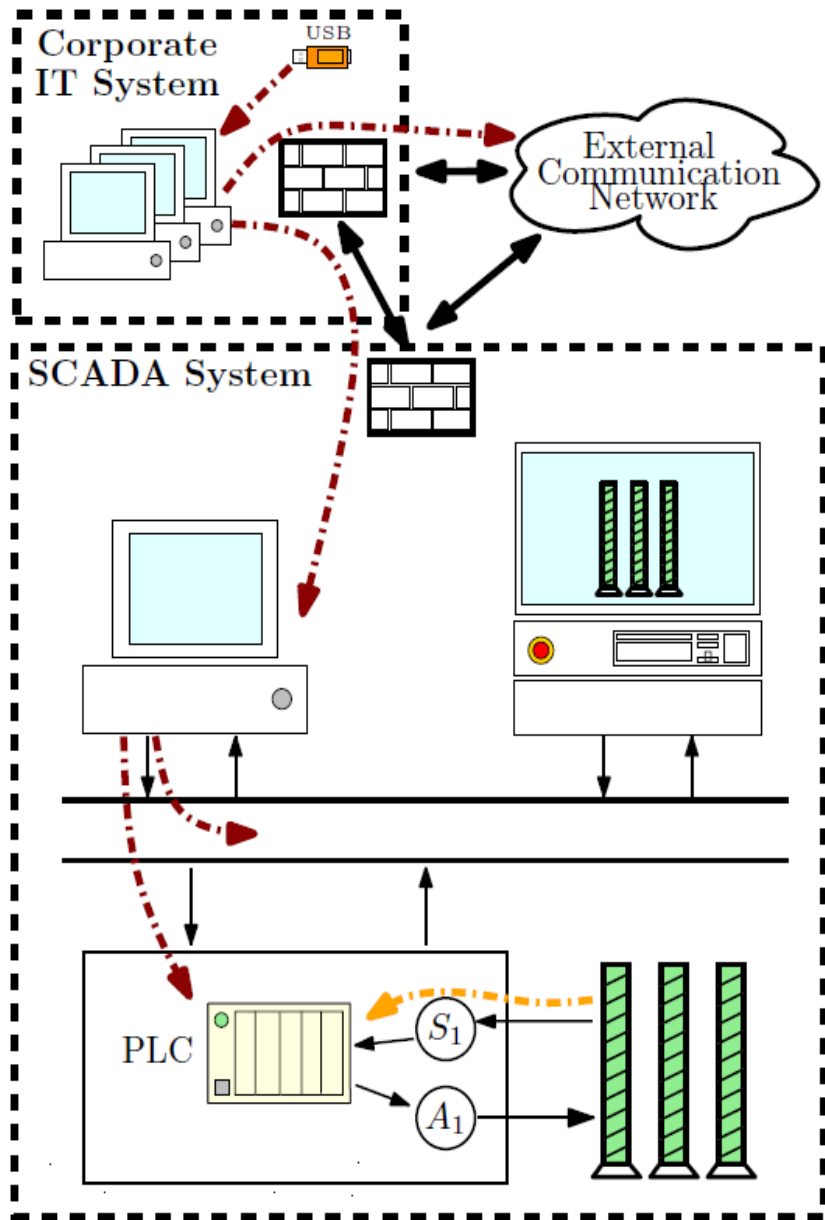
## Example 2: The Stuxnet Worm (2010)

**Targets:** Windows, ICS, and PLCs connected to variable-frequency drives  
Exploited **4 zero-day flaws**

- **Speculated goal:**  
Harm centrifuges at uranium enrichment facility in Iran
- **Attack mode:**
  1. Delivery with USB stick (**no internet connection necessary**)
  2. Replay measurements to control center and execute harmful controls



[“The Real Story of Stuxnet”, IEEE Spectrum, 2013]  
(See also <http://www.zerodayfilm.com/> )



(a) Infection and data recording.

# Example 3: Events in Ukraine (December, 2015)

## Analysis confirms coordinated hack attack caused Ukrainian power outage

BlackEnergy was key ingredient used to cause power outage to at least 80k customers.

by Dan Goodin - Jan 11, 2016 5:42am GMT

[Share](#) [Tweet](#) [Email](#) **33**

The people who carried out **last month's first known hacker-caused power outage** used highly destructive malware to gain a foothold into multiple regional distribution power companies in Ukraine and delay restoration efforts once electricity had been shut off, a newly published analysis confirms.

The malware, known as BlackEnergy, allowed the attackers to gain a foothold on the power company systems, said the report, which was published by a member of the SANS

### FURTHER READING



The report stresses there's no evidence BlackEnergy or its recently developed KillDisk component was the direct cause of the outage, which so far has been shown to affect about 80,000 customers. The analysis also cautioned that evidence showing some past BlackEnergy infections relied on booby-trapped Microsoft Office documents to spread are no indication such a vector was used in the recent Ukrainian power grid attacks. Still, this weekend's report leaves little doubt the blackout was the result of a highly coordinated hacker attack that relied on BlackEnergy as a key ingredient.





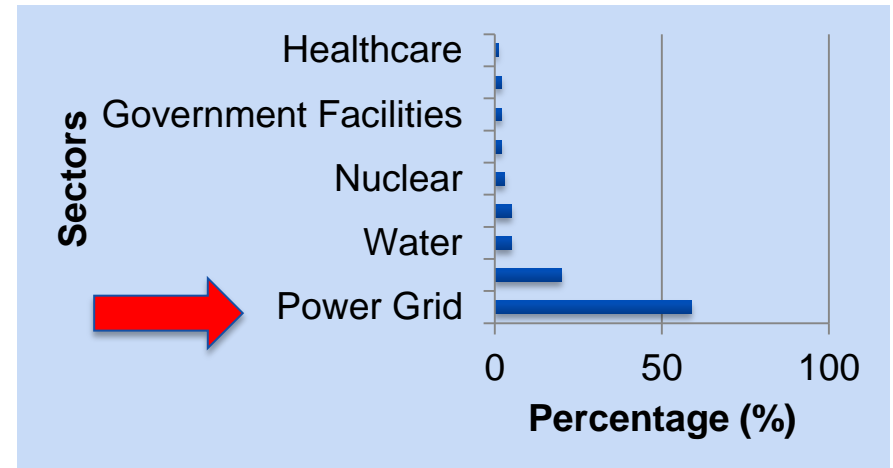
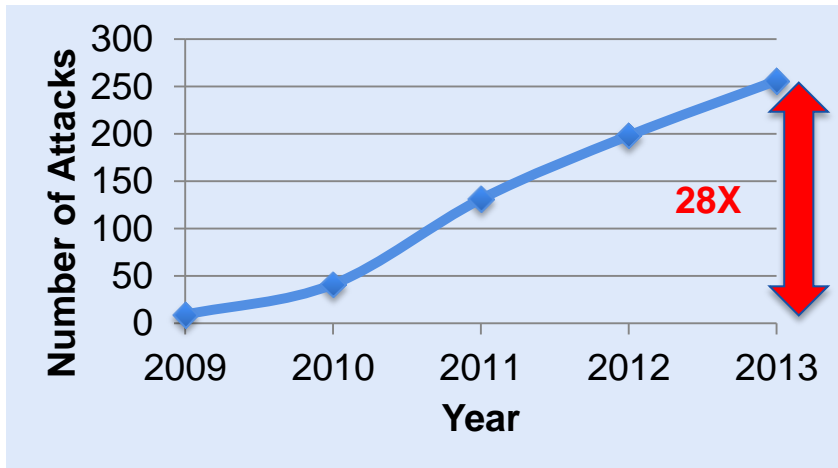
## Example 3: Events in Ukraine (December, 2015)

- BlackEnergy (2007-)
- From [arstechnica.com](http://arstechnica.com):
  - “In 2014 ... targeted the North Atlantic Treaty Organization, Ukrainian and Polish government agencies, and a variety of sensitive European industries”
  - “booby-trapped macro functions embedded in Microsoft Office documents”
  - “render infected computers unbootable”
  - “KillDisk, which destroys critical parts of a computer hard drive”
  - “backdoored secure shell (SSH) utility that gives attackers permanent access to infected computers”
- More advanced, more autonomous, follow-up attack in 2016: “Crash Override”



# Some Statistics

Cyber incidents in critical infrastructures in the US  
(Voluntarily reported to ICS-CERT)



[ICS-CERT, 2013]  
[S. Zonouz, 2014]

# Cyber-Physical Security

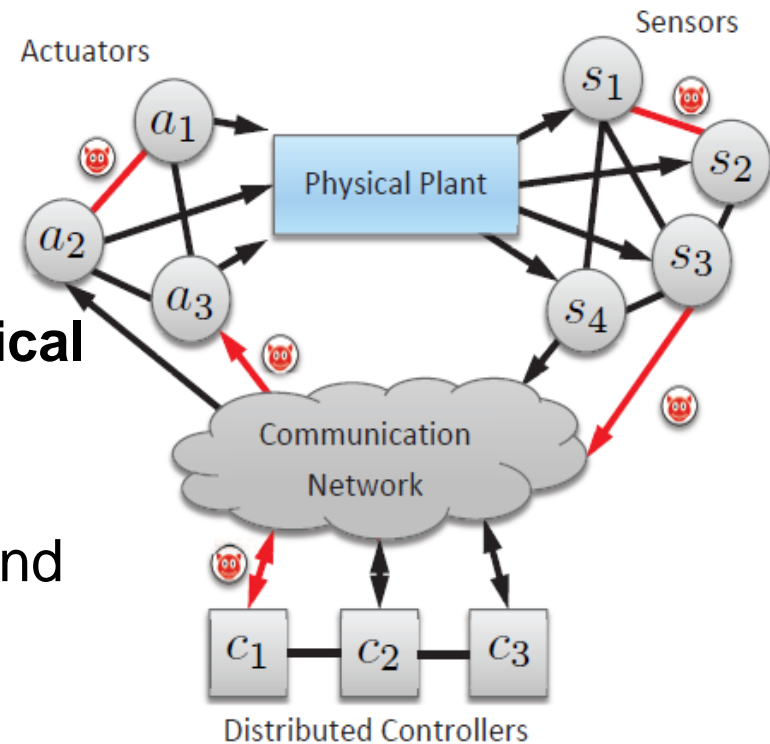
## Networked control systems

- are being **integrated with business/corporate networks**
- have many potential points of **cyber-physical attack**

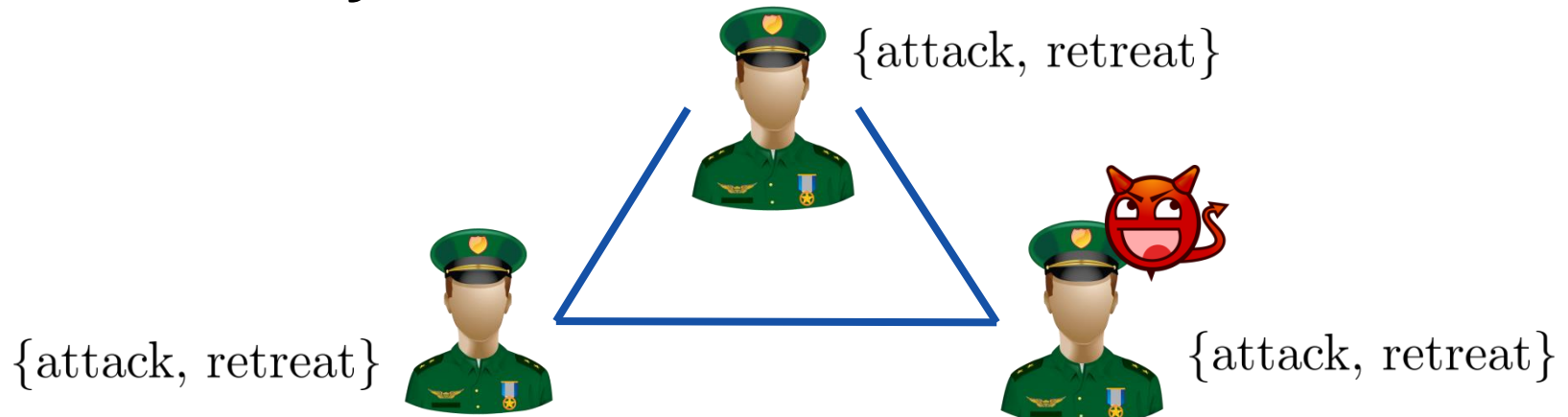
## Need tools and strategies to understand and mitigate attacks:

- Which threats should we care about?
- What impact can we expect from attacks?
- Which resources should we protect (more), and how?

## Is it enough to apply cyber (IT) security solutions?



# Example of Classic Cyber Security: The Byzantine Generals Problem



- Consider  $n$  generals and  $q$  unknown traitors among them. Can the  $n - q$  loyal generals always reach an agreement?
- Traitors (“Byzantine faults”) can do anything: different message to different generals, send no message, change forwarded message,...
- Agreement protocol exists iff  $n \geq 3q + 1$
- If loyal generals use unforgeable signed messages (“authentication”) then agreement protocol exists for any  $q$ ! [Lamport *et al.*, ACM TOPLAS, 1982]
- Application to linear consensus computations: See [Pasqualetti *et al.*, CDC, 2007], [Sundaram and Hadjicostis, ACC, 2008]

# Special Controls Perspective Needed?

Clearly cyber (IT) security is needed: Authentication, encryption, firewalls, etc.

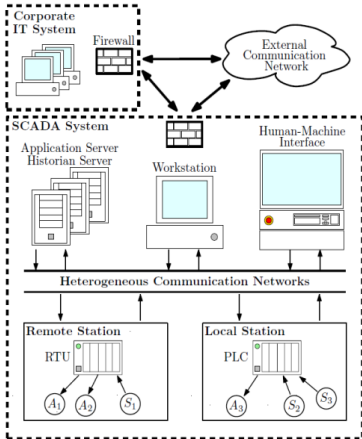
**But not sufficient...**

Interaction between physical and cyber systems make control systems different from normal IT systems

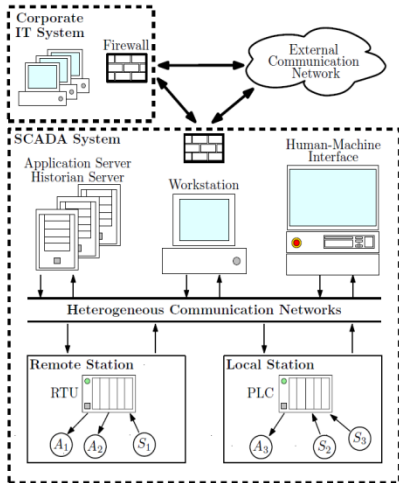
Malicious actions can enter anywhere in the closed loop and cause harm, whether channels secured or not

Can we trust the interfaces and channels are really secured?  
(see OpenSSL Heartbleed bug...)

[Cardenas *et al.*, 2008]



# Security Challenges in ICS



## “New” vulnerabilities and “new” threats:

- **Controllers are computers** (Relays → Microprocessors)
- **Networked** (Access from corporate network)
- **Commodity IT solutions** (Windows, TCP/IP,...)
- **Open design** (Protocols known)
- **Increasing size and functionality** (New services, wireless,...)
- **Large and highly skilled IT global workforce** (More IT knowledge)
- **Cybercrime** (Attack tools available)

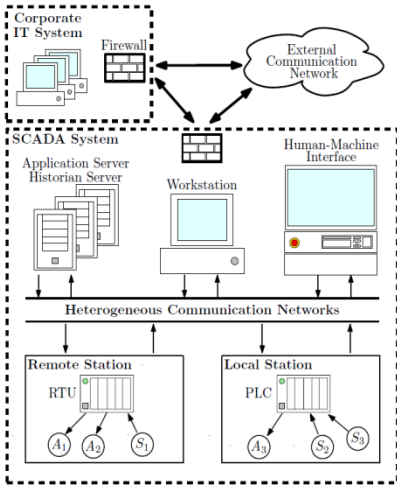
[Cardenas *et al.*, 2008]

# Security Challenges in ICS

## Differences to traditional IT systems:

- **Patching and frequent updates are not well suited for control systems**
- **Real-time availability** (Strict operational environment, sensitive to time delays. Designed for safety and easy access)
- **Legacy systems** (Often no authentication or encryption)
- **Protection of information and physical world** (Estimation and control algorithms)
- **Simpler network dynamics** (fixed topology, regular communication, limited number of protocols,...)

[Cardenas *et al.*, 2008]

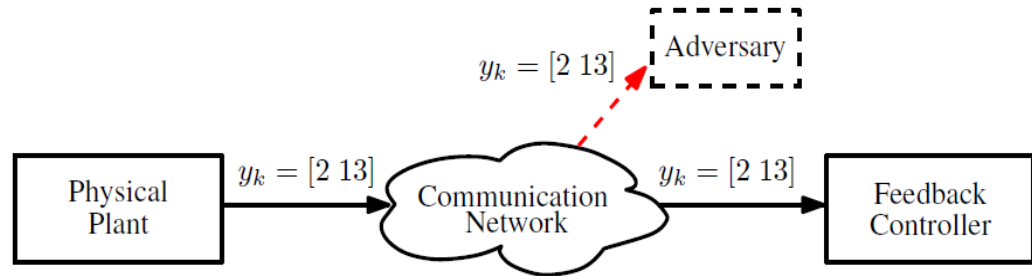


# CIA in Cyber Security [Bishop, 2002]

## C – Confidentiality

### “Privacy”

(See recent work by Le Ny, Pappas, Dullerud, Cortes, Tanaka, Sandberg,... )

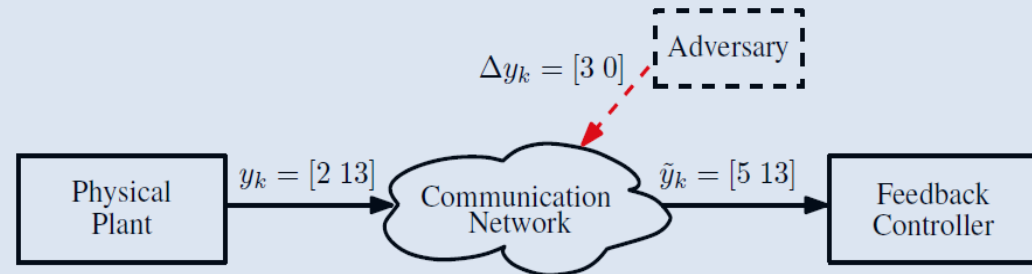


(a) Data confidentiality violation by a disclosure attack.

## I – Integrity

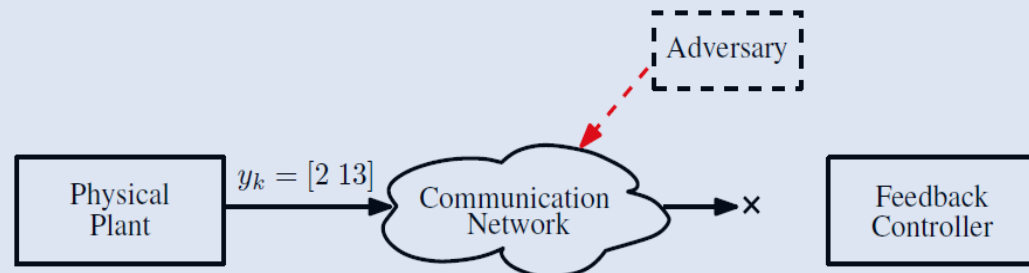
### “Security”

(Focus here. Good intro in CSM 2015 special issue)



(b) Data integrity violation by a false-data injection attack.

## A – Availability



(c) Data availability violation by a denial-of-service attack.



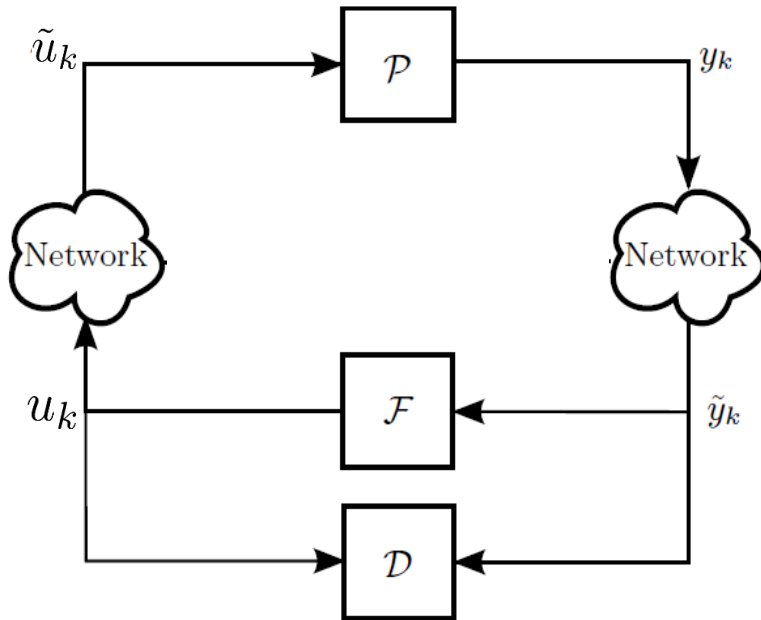


# Outline

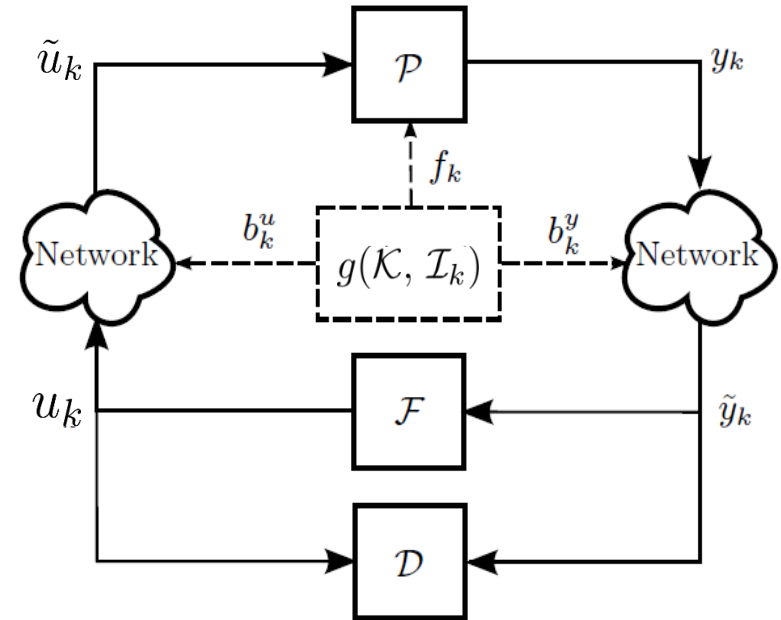
- Background and motivation
- **CPS attack models**
- Risk management
- Attack detectability and security metrics
- Attack identification and secure state estimation



# Networked Control System under Attack



- Physical plant ( $\mathcal{P}$ )
- Feedback controller ( $\mathcal{F}$ )
- Anomaly detector ( $\mathcal{D}$ )
- Disclosure Attacks

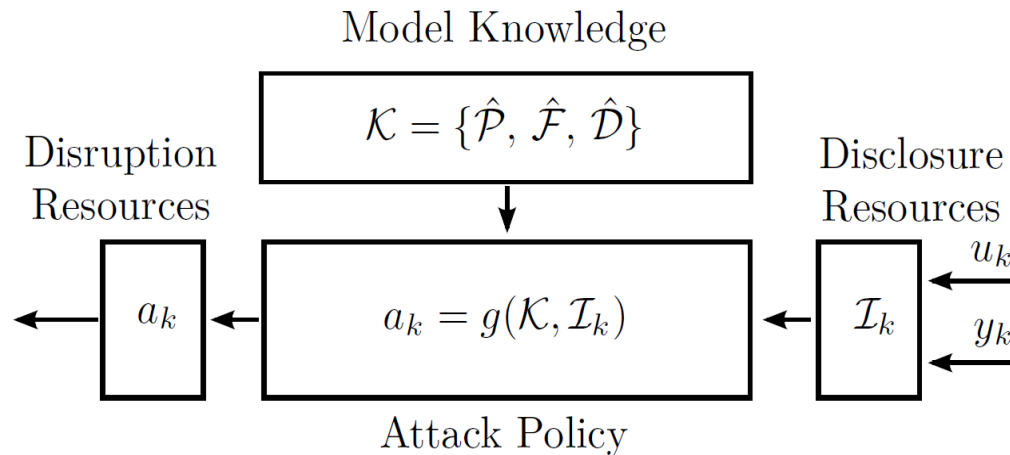


- Physical Attacks  $f_k$
- Deception Attacks

$$\tilde{u}_k = u_k + \Gamma^u b_k^u$$

$$\tilde{y}_k = y_k + \Gamma^y b_k^y$$

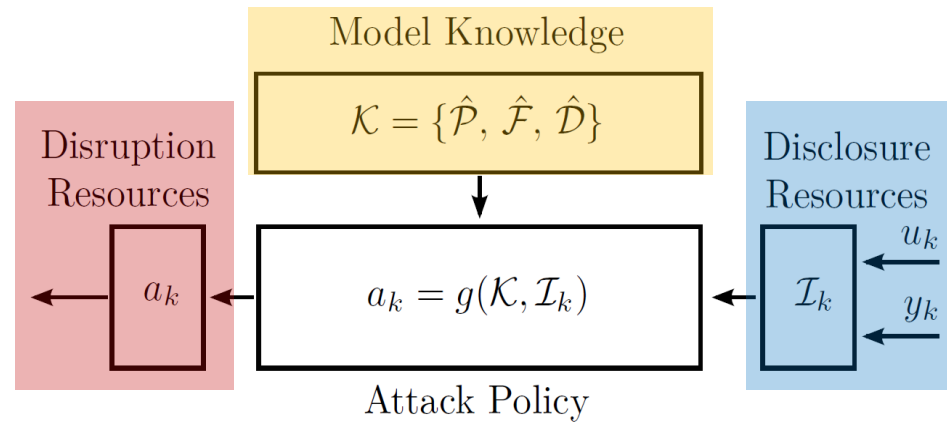
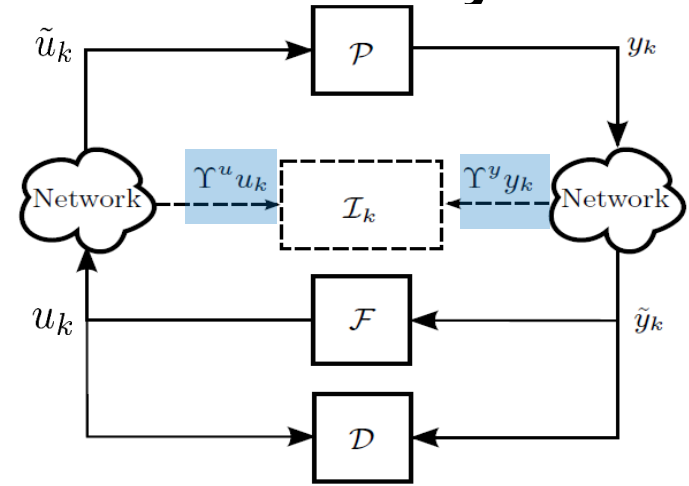
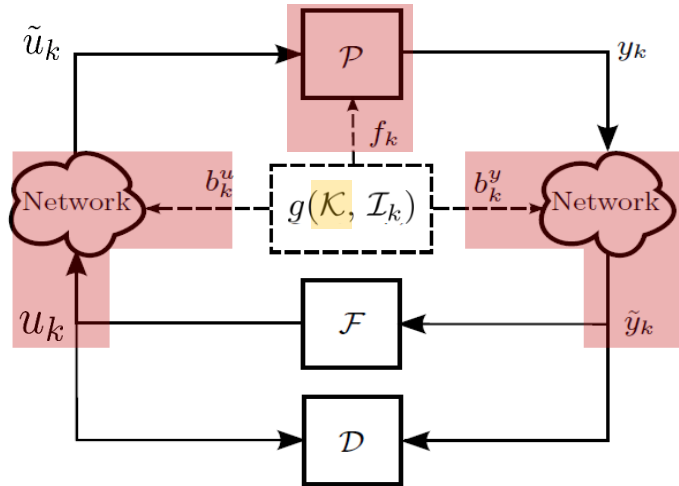
# Adversary Model



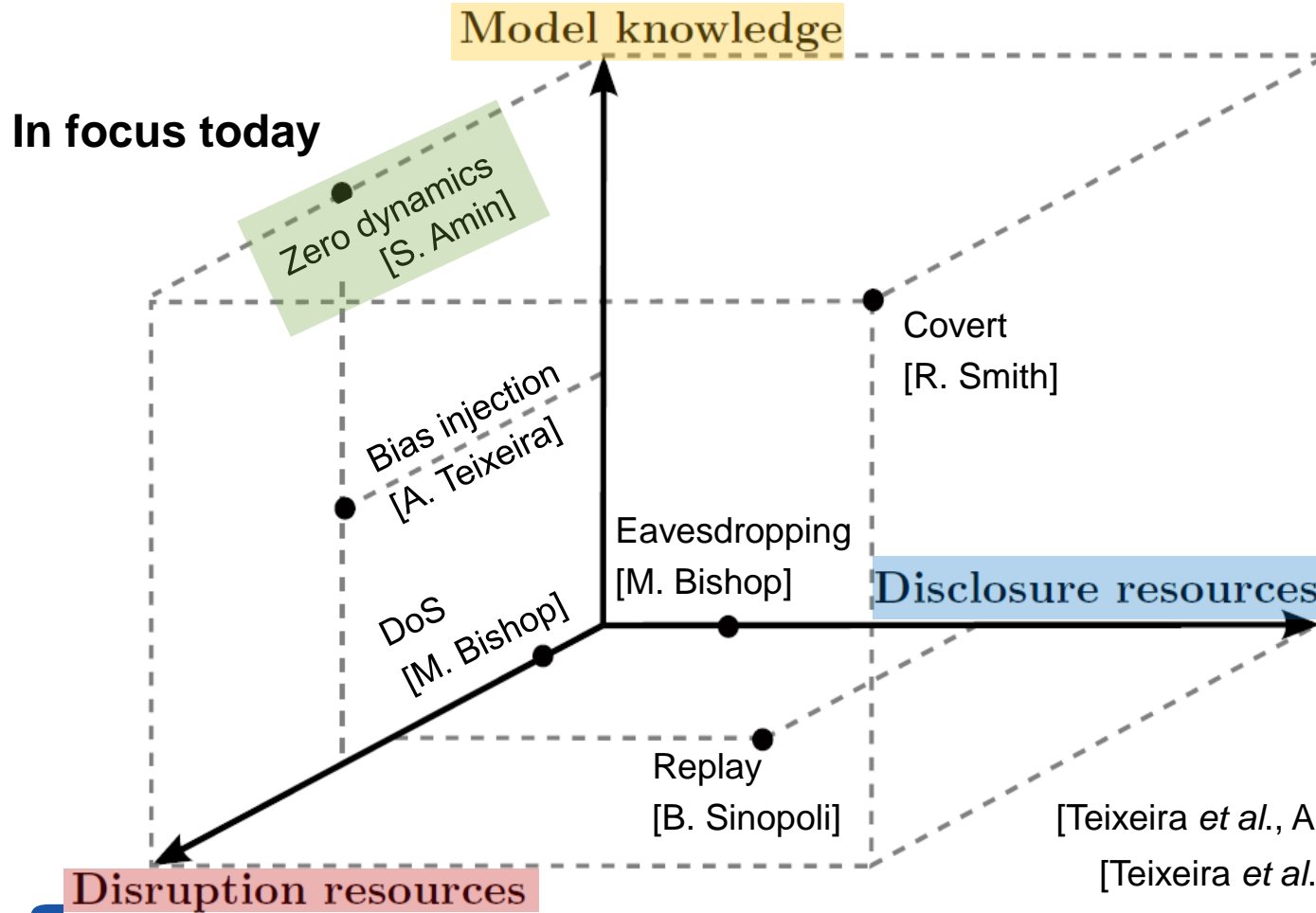
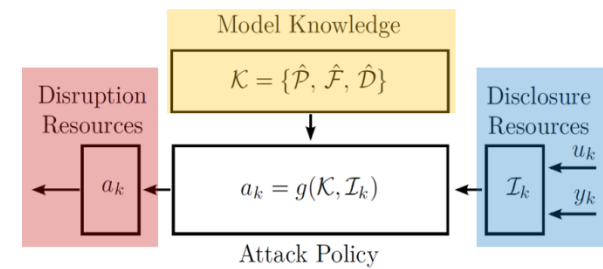
- **Attack policy:** Goal of the attack? Destroy equipment, increase costs,...
- **Model knowledge:** Adversary knows models of plant and controller? Possibility for stealthy attacks...
- **Disruption/disclosure resources:** Which channels can the adversary access?

[Teixeira *et al.*, HiCoNS, 2012]

# Networked Control System with Adversary Model



# Attack Space



[Teixeira *et al.*, Automatica, 2015]

[Teixeira *et al.*, HiCoNS, 2012]



# Outline

- Background and motivation
- CPS attack models
- **Risk management**
- Attack detectability and security metrics
- Attack identification and secure state estimation

# Why Risk Management?

Complex control systems with numerous attack scenarios

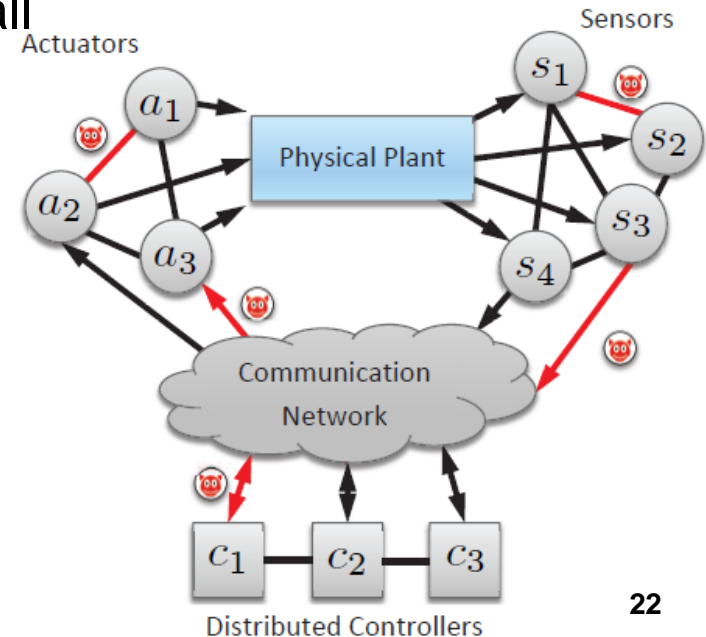
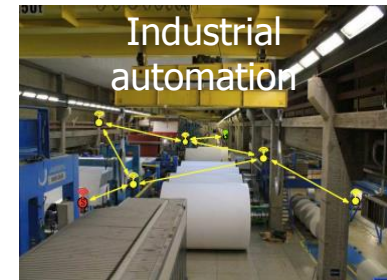
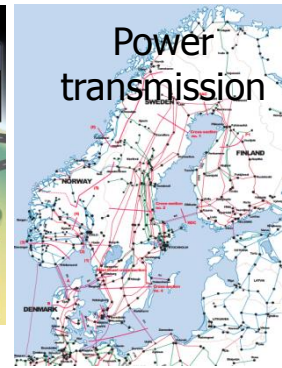
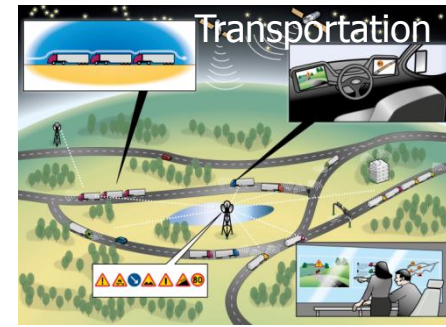
Examples: Critical infrastructures (power, transport, water, gas, oil) often with weak security guarantees

Too costly to secure the entire system against all attack scenarios

What scenarios to prioritize?

What components to protect?

When possible to identify attacks?





# Defining Risk

**Risk = (Scenario, Likelihood, Impact)**

## Scenario

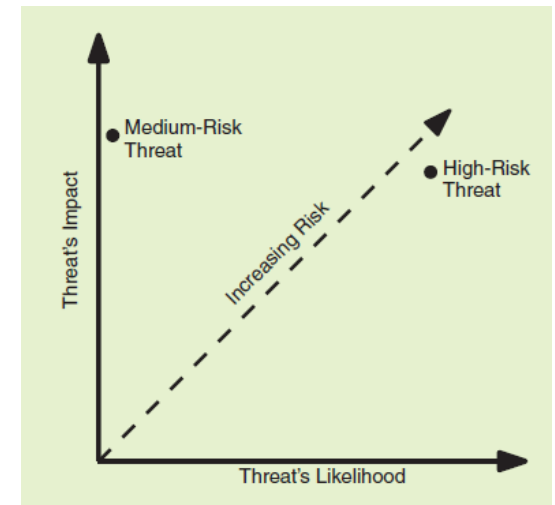
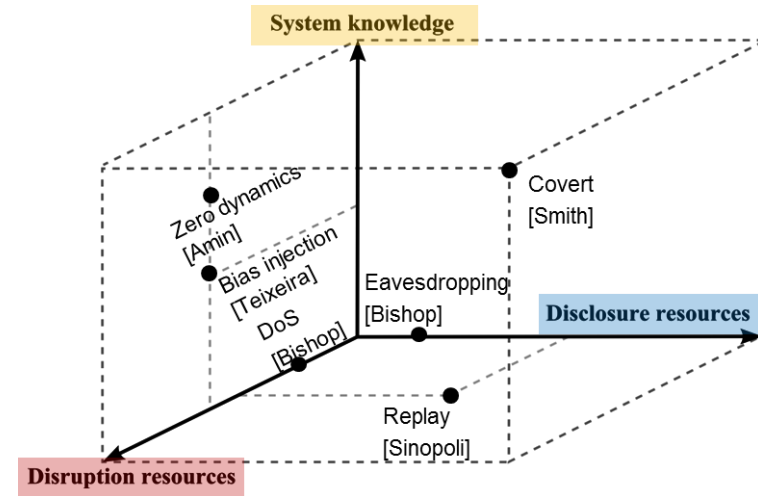
- How to describe the system under attack?

## Likelihood

- How much effort does a given attack require?

## Impact

- What are the consequences of an attack?

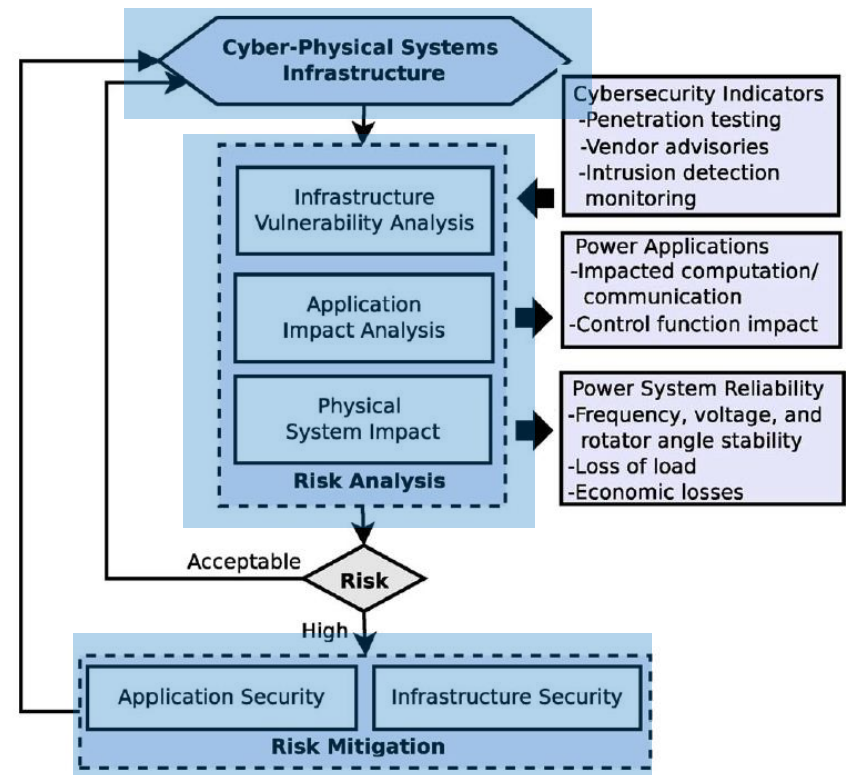


[Kaplan & Garrick, 1981], [Bishop, 2002]  
 ([Teixeira *et al.*, IEEE CSM, 2015])

# Risk Management Cycle

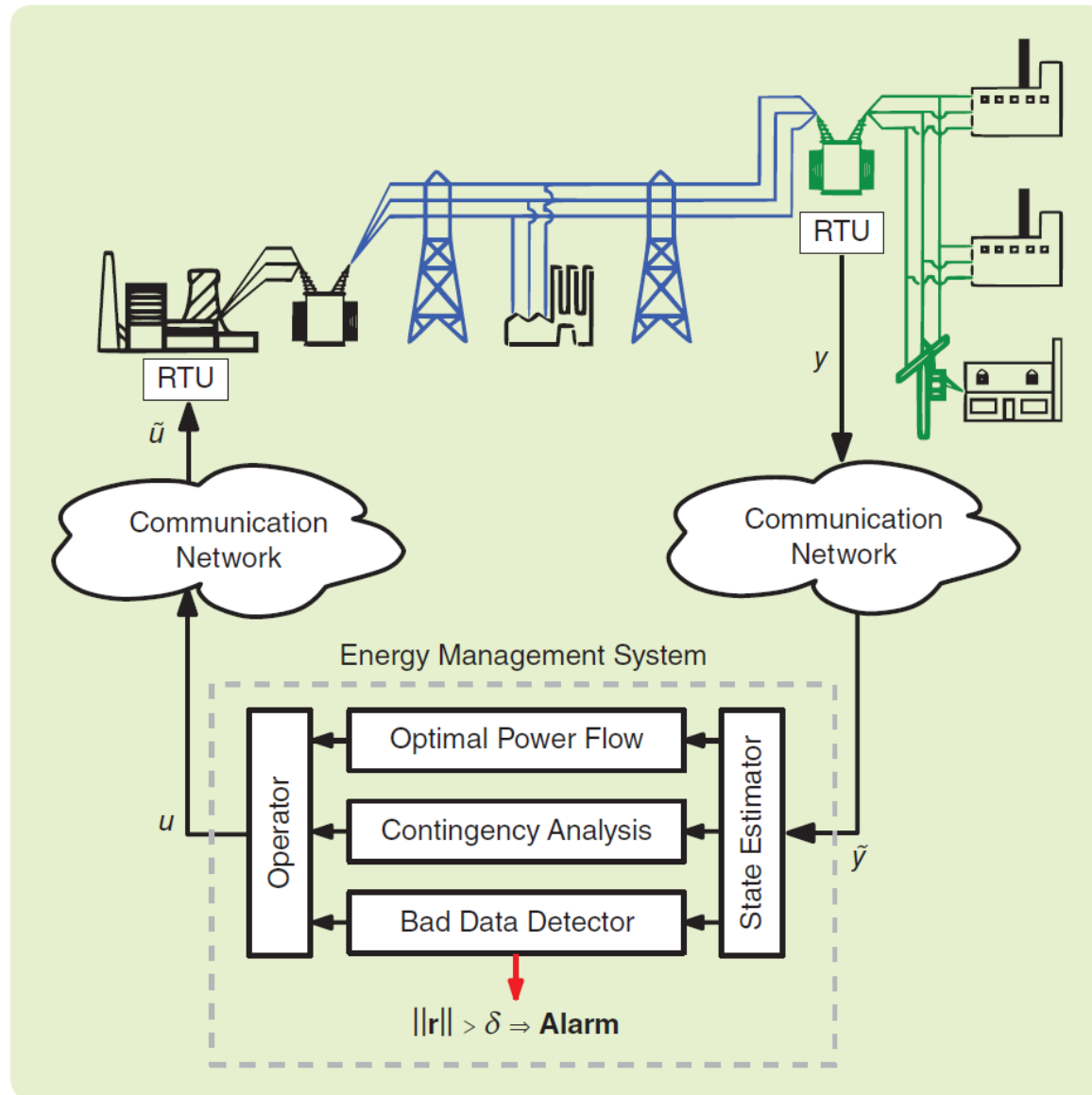
## Main steps in risk management

- Scope definition
  - Models, Scenarios, Objectives
- Risk Analysis
  - **Threat Identification**
  - **Likelihood Assessment**
  - Impact Assessment
- Risk Treatment
  - **Prevention, Detection, Mitigation**

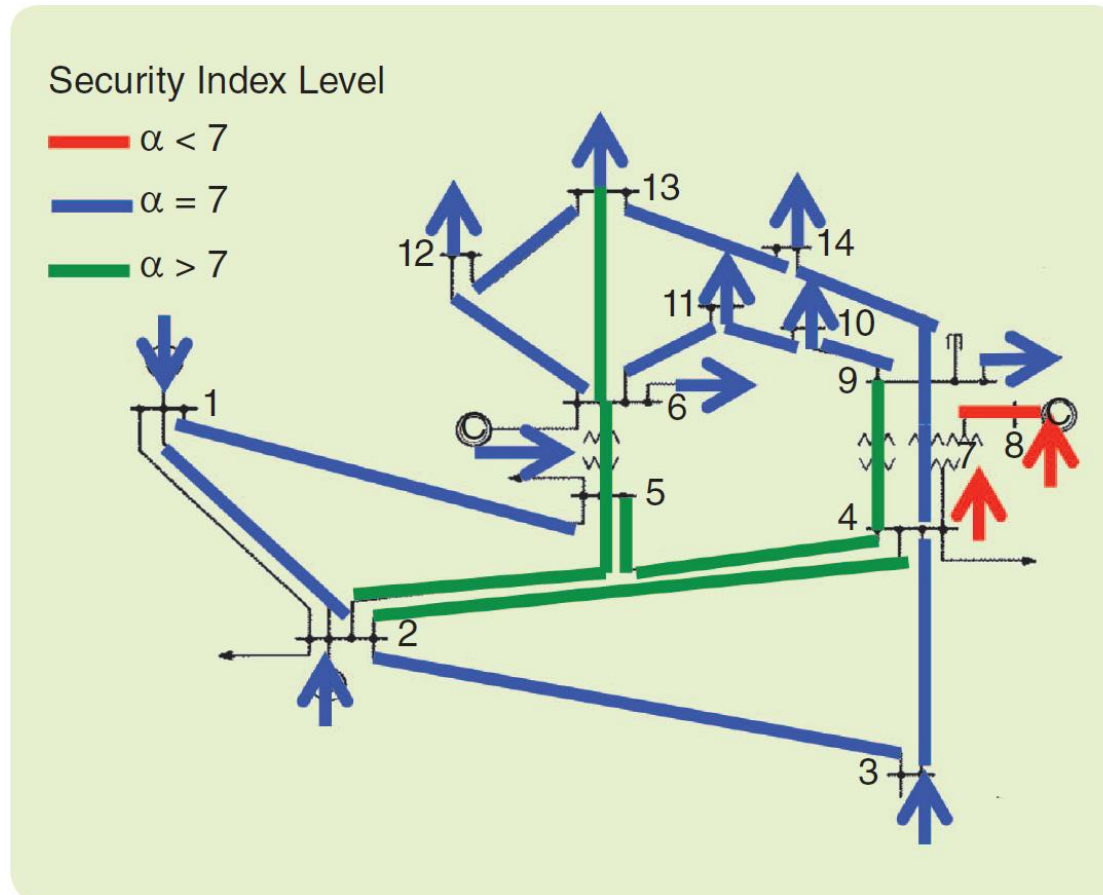


[Sridhar *et al.*, Proc. IEEE, 2012]

# Example: Power System State Estimator



# Example: Power System State Estimator



Security index  $\alpha$  (to be defined) indicates sensors with inherent weak redundancy ( $\sim$ security). These should be secured first!



# Outline

- Background and motivation
- CPS attack models
- Risk management
- **Attack detectability and security metrics**
- Attack identification and secure state estimation



# Basic Notions: Input Observability and Detectability

$$\begin{aligned}x(k+1) &= Ax(k) + Bu(k), & x(k) &\in \mathbb{R}^n, u(k) \in \mathbb{R}^m \\y(k) &= Cx(k) + Du(k), & y(k) &\in \mathbb{R}^p\end{aligned}$$

## Definitions:

1. The input  $u$  is *observable with knowledge of  $x(0)$*  if  $y(k) = 0$  for  $k \geq 0$  implies  $u(k) = 0$  for  $k \geq 0$ , provided  $x(0) = 0$
2. The input  $u$  is *observable* if  $y(k) = 0$  for  $k \geq 0$  implies  $u(k) = 0$  for  $k \geq 0$  ( $x(0)$  unknown)
3. The input  $u$  is *detectable* if  $y(k) = 0$  for  $k \geq 0$  implies  $u(k) \rightarrow 0$  for  $k \rightarrow \infty$  ( $x(0)$  unknown)

[Hou and Patton, Automatica, 1998]



# Basic Notions: Input Observability and Detectability

The Rosenbrock system matrix:

$$P(z) = \begin{bmatrix} A - zI & B \\ C & D \end{bmatrix} \in \mathbb{C}^{(n+p) \times (n+m)}$$

## First observations:

- Necessary condition for Definitions 1-3  
 $\max_z \text{rank } P(z) = m + n \Leftrightarrow \text{normalrank } P(z) = m + n$
- Fails if number of inputs larger than number of outputs ( $m > p$ )
- Necessary and sufficient conditions involve the *invariant zeros*:  
 $\sigma(P(z)) := \{z : \text{rank } P(z) < \text{normalrank } P(z)\}$   
(Transmission zeros + uncontrollable/unobservable modes,  
Matlab command: `tzero`)





# Basic Notions: Input Observability and Detectability

$$P(z) = \begin{bmatrix} A - zI & B \\ C & D \end{bmatrix} \in \mathbb{C}^{(n+p) \times (n+m)}$$

**Theorems.** Suppose  $(A, B, C, D)$  is minimal realization.

1. The input  $u$  is *observable with knowledge of  $x(0)$*   $\Leftrightarrow$   
 $\max_z \text{rank } P(z) = m + n \Leftrightarrow \text{normalrank } P(z) = m + n$
2. The input  $u$  is *observable*  $\Leftrightarrow$   
 $\forall z : \text{rank } P(z) = m + n$   
(no invariant zeros)
3. The input  $u$  is *detectable*  $\Leftrightarrow$  (1) and  
 $\sigma(P(z)) \subseteq \{z : |z| < 1\}$   
(invariant zeros are all stable = system is minimum phase)

[Hou and Patton, Automatica, 1998]



# Basic Notions: Input Observability and Detectability

$$P(z) = \begin{bmatrix} A - zI & B \\ C & D \end{bmatrix}, \quad O(z) = \begin{bmatrix} A - zI \\ C \end{bmatrix}$$

**Theorems.**  $(A, B, C, D)$  possibly non-minimal realization

1. The input  $u$  is *observable with knowledge of  $x(0)$*   $\Leftrightarrow$   
 $\max_z \text{rank } P(z) = m + n \Leftrightarrow \text{normalrank } P(z) = m + n$
- 2'. The input  $u$  is *observable*  $\Leftrightarrow$  (1) and  
 $\sigma(P(z)) = \sigma(O(z))$   
(invariant zeros are all unobservable modes)
- 3'. The input  $u$  is *detectable*  $\Leftrightarrow$  (1) and  
 $\sigma(P(z)) \setminus \sigma(O(z)) \subseteq \{z : |z| < 1\}$   
(invariant zeros that are not unobservable modes are all stable)

[Hou and Patton, Automatica, 1998]



# Fault Detection vs. Secure Control

## Typical condition used in fault detection/fault tolerant control:

1. The input  $u$  is *observable with knowledge of  $x(0)$*   $\Leftrightarrow$   
$$\max_z \text{rank } P(z) = m + n \Leftrightarrow \text{normalrank } P(z) = m + n$$

[Ding, Patton]

## Typical conditions used in secure control/estimation:

2. The input  $u$  is *observable*  $\Leftrightarrow$   
$$\forall z : \text{rank } P(z) = m + n$$
  
(no invariant zeros)

[Sundaram, Tabuada]

- 3/3'. The input  $u$  is *detectable*  $\Leftrightarrow$  (1) and

[Pasqualetti, Sandberg]

$$\sigma(P(z)) \subseteq \{z : |z| < 1\}$$

(invariant zeros are all stable = system is minimum phase)

## Example

$$A = \begin{pmatrix} 0.9 & 0 & 0 \\ 0 & 0.8 & 0 \\ 0 & 0 & 0.9 \end{pmatrix}, B = \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \\ 0 & 0.25 \end{pmatrix}, C = \begin{pmatrix} 0.4 & 0.6 & 0 \\ 0.2 & 0 & 0.4 \end{pmatrix}$$

$$G(z) = C(zI - A)^{-1}B + D = \begin{pmatrix} \frac{0.2}{z-0.9} & \frac{0.3}{z-0.8} \\ \frac{0.1}{z-0.9} & \frac{0.1}{z-0.9} \end{pmatrix}$$

Invariant zeros =  $\sigma(P(z)) = \{1.1\}$

1. The input  $u$  is *observable with knowledge of  $x(0)$* : **Yes!**
2. The input  $u$  is *observable*: **No!**
3. The input  $u$  is *detectable*: **No!**

With  $x(0) = \begin{pmatrix} -0.705 \\ 0.470 \\ 0.352 \end{pmatrix}$  and  $u(k) = 1.1^k \begin{pmatrix} -0.282 \\ 0.282 \end{pmatrix}$  then  $y(k) = 0, k \geq 0$

**OK for fault detection but perhaps not for security!**

# Attack and Disturbance Model

Consider the linear system  $y = G_d d + G_a a$  (the controlled infrastructure):

$$x(k+1) = Ax(k) + B_d d(k) + B_a a(k)$$

$$y(k) = Cx(k) + D_d d(k) + D_a a(k)$$

- Unknown state  $x(k) \in \mathbb{R}^n$  ( $x(0)$  in particular)
- Unknown (natural) disturbance  $d(k) \in \mathbb{R}^o$
- Unknown (malicious) attack  $a(k) \in \mathbb{R}^m$
- Known measurement  $y(k) \in \mathbb{R}^p$
- Known model  $A, B_d, B_a, C, D_d, D_a$
  
- **Definition:** Attack signal  $a$  is *persistent* if  $a(k) \not\rightarrow 0$  as  $k \rightarrow \infty$
  
- **Definition:** A (persistent) attack signal  $a$  is *undetectable* if there exists a simultaneous (masking) disturbance signal  $d$  and initial state  $x(0)$  such that  $y(k) = 0, k \geq 0$  (Cf. Theorem 3')

# Undetectable Attacks and Masking

The Rosenbrock system matrix:

$$P(z) = \begin{bmatrix} A - zI & B_d & B_a \\ C & D_d & D_a \end{bmatrix}$$

- Attack signal  $a(k) = z_0^k a_0$ ,  $0 \neq a_0 \in \mathbb{C}^m$ ,  $z_0 \in \mathbb{C}$ , is *undetectable* iff there exists  $x_0 \in \mathbb{C}^n$  and  $d_0 \in \mathbb{C}^o$  such that

$$P(z_0) \begin{bmatrix} x_0 \\ d_0 \\ a_0 \end{bmatrix} = 0$$

- Attack signal is undetectable if indistinguishable from measurable ( $y$ ) effects of natural noise ( $d$ ) or uncertain initial states ( $x_0$ ) [**masking**]

## Example (cont'd)

$$G(z) = C(zI - A)^{-1}B + D = (G_d(z) \quad G_a(z))$$

$$G_d(z) = \begin{pmatrix} \frac{0.2}{z-0.9} \\ \frac{0.1}{z-0.9} \end{pmatrix}, \quad G_a(z) = \begin{pmatrix} \frac{0.3}{z-0.8} \\ \frac{0.1}{z-0.9} \end{pmatrix}$$

Poles =  $\{0.9, 0.9, 0.8\}$

Invariant zeros =  $\sigma(P(z)) = \{1.1\}$

Undetectable attack:  $a(k) = 1.1^k \cdot 0.282$

Masking initial state:  $x_0 = \begin{pmatrix} -0.705 \\ 0.470 \\ 0.352 \end{pmatrix}$

Masking disturbance  $d(k) = 1.1^k \cdot (-0.282)$





# Undetectable Attacks and Masking (cont'd)

- Suppose operator observes the output  $y(k)$ , and does *not know* the true initial state  $x(0)$  and true disturbance  $d(k)$
- Let  $(x_0, d_0, a_0)$  be an undetectable attack,  $0 = G_d d_0 + G_a a_0$  with initial state  $x_0$

Consider the cases:

1. **Un-attacked system**  $y = G_d(-d_0)$ , with initial state  $x(0) = 0$
2. **Attacked system**  $y = G_a a_0$ , with initial state  $x(0) = x_0$

If initial states  $x(0) = 0$  and  $x(0) = x_0$  and disturbances  $d = -d_0$  and  $d = 0$  are equally likely, then impossible for operator to decide which case is true  $\Rightarrow$  **Attack is undetectable!**



# Undetectable Attacks and Masking (cont'd)

- Suppose operator observes the output  $y(k)$ , and does *not know* the true initial state  $x(0)$  and true disturbance  $d(k)$
- Let  $(x_0, d_0, a_0)$  be an undetectable attack,  $0 = G_d d_0 + G_a a_0$  with initial state  $x_0$

Consider the cases:

1. **Un-attacked system 1:**  $y = G_d d$ , with initial state  $x(0)$
2. **Attacked system:**  $y = G_d(d + d_0) + G_a a_0$ , with initial state  $x(0) + x_0$

If initial states  $x(0)$  and  $x(0) + x_0$  and disturbances  $d$  and  $d + d_0$  are equally likely, then impossible for operator to decide which case is true  $\Rightarrow$  **Attack is undetectable!**



# Undetectable Attacks and Masking (cont'd)

- Suppose operator observes the output  $y(k)$ , and does *not know* the true initial state  $x(0)$  and true disturbance  $d(k)$
- Let  $(x_0, d_0, a_0)$  be an undetectable attack,  $0 = G_d d_0 + G_a a_0$  with initial state  $x_0$

Consider the cases:

1. **Un-attacked system**  $y = G_d d$ , with initial state  $x(0)$
2. **Attacked system**  $y = G_d(d + d_0) + G_a a_0$ , with initial state  $x(0) + x_0$

If initial states  $x(0)$  and  $x(0) + x_0$  and disturbances  $d$  and  $d + d_0$  are equally likely, then impossible for operator to decide which case is true  $\Rightarrow$  **Attack is undetectable!**

# Undetectable Attacks and Masking (cont'd)

- Suppose operator observes the output  $y(k)$ , and *knows* the true initial state  $x(0) = 0$  and the disturbance  $d(k) = 0, k \geq 0$
- Suppose system is asymptotically stable,  $\rho(A) < 1$
- Let  $(x_0, a_0)$  be an undetectable attack,  $0 = G_a a_0$  with initial state  $x_0$

Consider the cases:

1. **Un-attacked system**  $y_1(k) = 0, k \geq 0$ , with initial state  $x(0) = 0$
2. **Attacked system**  $y_2(k) = (G_a a_0)(k) = -CA^k x_0 \rightarrow 0$  as  $k \rightarrow \infty$ , with initial state  $x(0) = 0$

The attacked output  $y_2$  is vanishing, and can be made arbitrarily close to  $y_1$  by scaling  $(x_0, a_0) \Rightarrow$  **Attack is asymptotically undetectable!**

# The Security Index $\alpha_i$

$$\alpha_i := \min_{|z_0| \geq 1, x_0, d_0, a_0^i} \|a_0^i\|_0$$

subject to  $P(z_0) \begin{bmatrix} x_0 \\ d_0 \\ a_0^i \end{bmatrix} = 0$

**Notation:**  $\|a\|_0 := |\text{supp}(a)|$ ,  $a^i$  vector  $a$  with  $i$ -th element non-zero

## Interpretation:

- Attacker persistently targets signal component  $a_i$  (condition  $|z_0| \geq 1$ )
- $\alpha_i$  is smallest number of attack signals that need to be simultaneously accessed to stage undetectable attack against signal  $a_i$

**Argument:** Large  $\alpha_i \Rightarrow$  malicious cyber attacks targeting  $a_i$  less likely

Problem NP-hard in general (combinatorial optimization, cf. matrix *spark*).  
Generalization of static index in [Sandberg *et al.*, SCS, 2010]

## Simple Example of Security Index

$$P(z) = \begin{bmatrix} A - zI & B_d & B_a \\ 0 & 0 & D_a \end{bmatrix} \quad D_a = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

- Measurements not affected by physical states and disturbances
- 3 measurements
- 4 attacks with security indices:
  - $\alpha_1 = 3$
  - $\alpha_2 = 3$
  - $\alpha_3 = 3$
  - $\alpha_4 = \infty$  (By definition. Even access to all attack signals not enough to hide attack)

# Special Case 1: Critical Attack Signals

Signal with  $\alpha_i = 1$  can be undetectably attacked without access to other elements  $\Rightarrow$  **Critical Attack Signal**

$$P_i(z) = \begin{bmatrix} A - zI & B_d & B_{a,i} \\ C & D_d & D_{a,i} \end{bmatrix} \in \mathbb{C}^{(n+p) \times (n+o+1)}, \quad P_d(z) = \begin{bmatrix} A - zI & B_d \\ C & D_d \end{bmatrix} \in \mathbb{C}^{(n+p) \times (n+o)}$$

**Simple test,  $\forall i$ :** If there is  $z_0 \in \mathbb{C}$ ,  $|z_0| \geq 1$ , such that  $\text{rank} [P_d(z_0)] = \text{rank} [P_i(z_0)]$ , then  $\alpha_i = 1$

**Even more critical case:** If  $\text{normalrank} [P_d(z_0)] = \text{normalrank} [P_i(z_0)]$  then there is undetectable critical attack for all frequencies  $z_0$

Holds generically when more disturbances than measurements ( $o \geq p$ )!

**Protect against these attack signals first in risk management!**

## Special Case 2: Transmission Zeros

$$P(z) = \begin{bmatrix} A - zI & B_d & B_a \\ C & D_d & D_a \end{bmatrix} \quad \begin{array}{l} \text{[Amin et al., ACM HSCC, 2010]} \\ \text{[Pasqualetti et al., IEEE TAC, 2013]} \end{array}$$

Suppose  $P(z)$  has full column normal rank. Then undetected attacks only at finite set of transmission zeros  $\{z_0\}$

$$\begin{array}{l} \text{Solve} \quad \alpha_i := \min_{|z_0| \geq 1, x_0, d_0, a_0^i} \|a_0^i\|_0 \\ \text{subject to} \quad P(z_0) \begin{bmatrix} x_0 \\ d_0 \\ a_0^i \end{bmatrix} = 0 \end{array}$$

by inspection of corresponding zero directions  $\Rightarrow$  **Easy in typical case of 1-dimensional zero directions**



## Special Case 3: Sensor Attacks

$$P(z) = \begin{bmatrix} A - zI & 0 & 0 \\ C & D_d & D_a \end{bmatrix}$$

[Fawzi *et al.*, IEEE TAC, 2014]  
[Chen *et al.*, IEEE ICASSP, 2015]  
[Lee *et al.*, ECC, 2015]

$P(z)$  only loses rank in eigenvalues  $z_0 \in \{\lambda_1(A), \dots, \lambda_n(A)\}$

Simple eigenvalues give one-dimensional spaces of eigenvectors  $x_0 \Rightarrow$  **Simplifies computation of  $\alpha_i$**

**Example:** Suppose  $D_a = I_p$  (sensor attacks),  $D_d = 0$ , and system observable from each  $y_i, i = 1, \dots, p$ :

- By the PBH-test:  $\alpha_i = p$  or  $\alpha_i = +\infty$  (if all eigenvalues stable, no persistent undetectable sensor attack exists)
- Redundant measurements increase  $\alpha_i$ !



# Special Case 4: Sensor Attacks for Static Systems

$$P(z) = \begin{bmatrix} I - zI & 0 & 0 \\ C & 0 & D_a \end{bmatrix} \quad \begin{array}{l} \text{[Liu et al., ACM CCS, 2009]} \\ \text{[Sandberg et al., SCS, 2010]} \end{array}$$

Since  $A = I_n$  and  $B_d = B_a = 0$ , this is the steady-state case

Space of eigenvectors  $x_0$  is  $n$ -dimensional  $\Rightarrow$  **Typically makes computation of  $\alpha_i$  harder than in the dynamical case!**

Practically relevant case in power systems where  $p > n \gg 0$

- Problem NP-hard, but power system imposes special structures in  $C$  (unimodularity etc.)
- Several works on efficient and exact computation of  $\alpha_i$  using min-cut/max-flow and  $\ell_1$ -relaxation ([Hendrickx et al., 2014], [Kosut, 2014], [Yamaguchi et al., 2015])



# Special Case 4: Solution by MILP

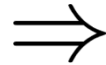
Big  $M$  reformulation:

$$\alpha_i := \min_{x_0, a_0} \|a_0\|_0$$

subject to

$$0 = Cx_0 + D_a a_0$$

$$a_{0,i} = 1$$



$$\alpha_i := \min_{z_0, x_0, a_0, z_k} \sum_k z_k$$

subject to

$$0 = Cx_0 + D_a a_0$$

$$a_{0,i} = 1$$

$$-Mz \leq a_0 \leq Mz$$

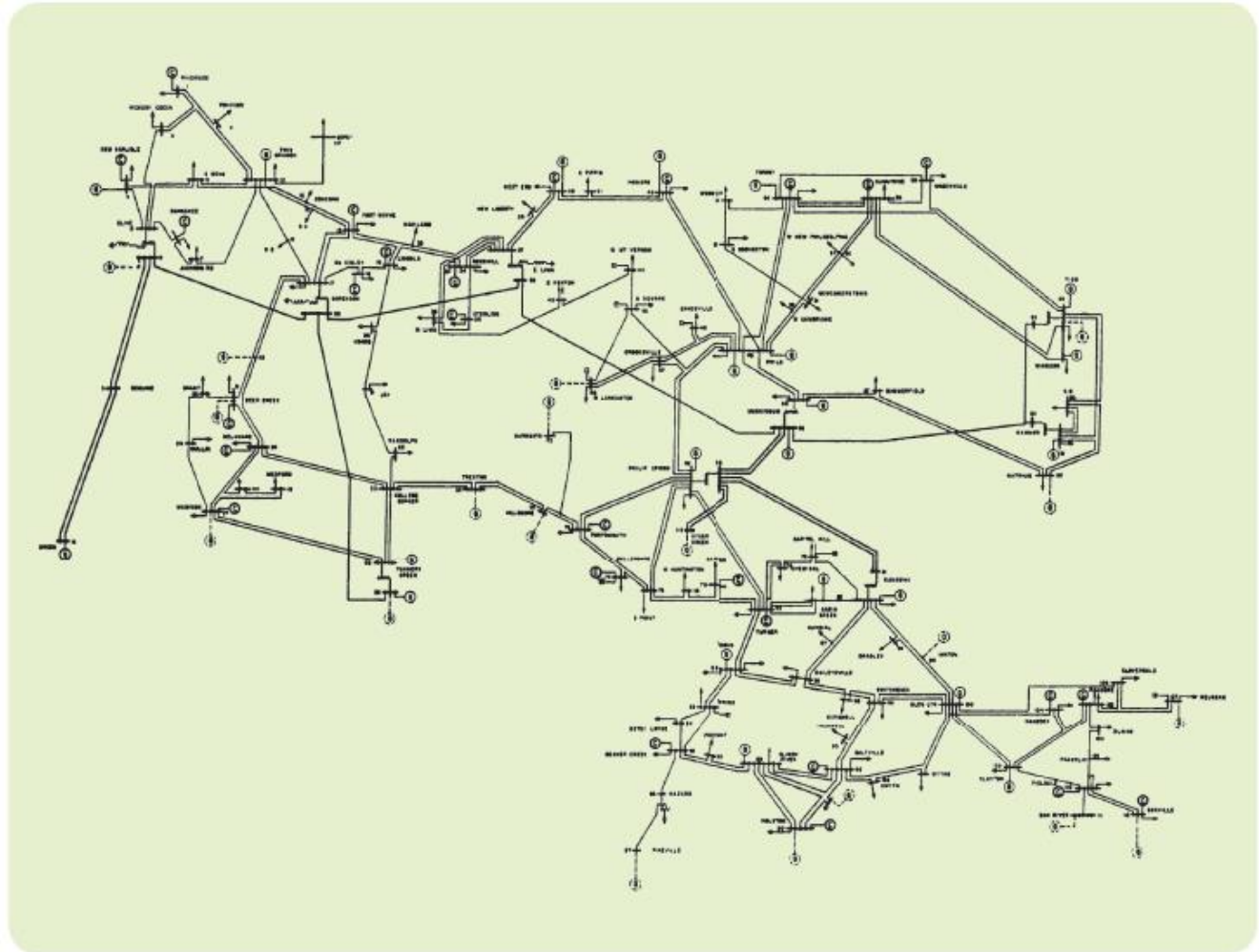
$$z_k \in \{0, 1\}$$

Elementwise

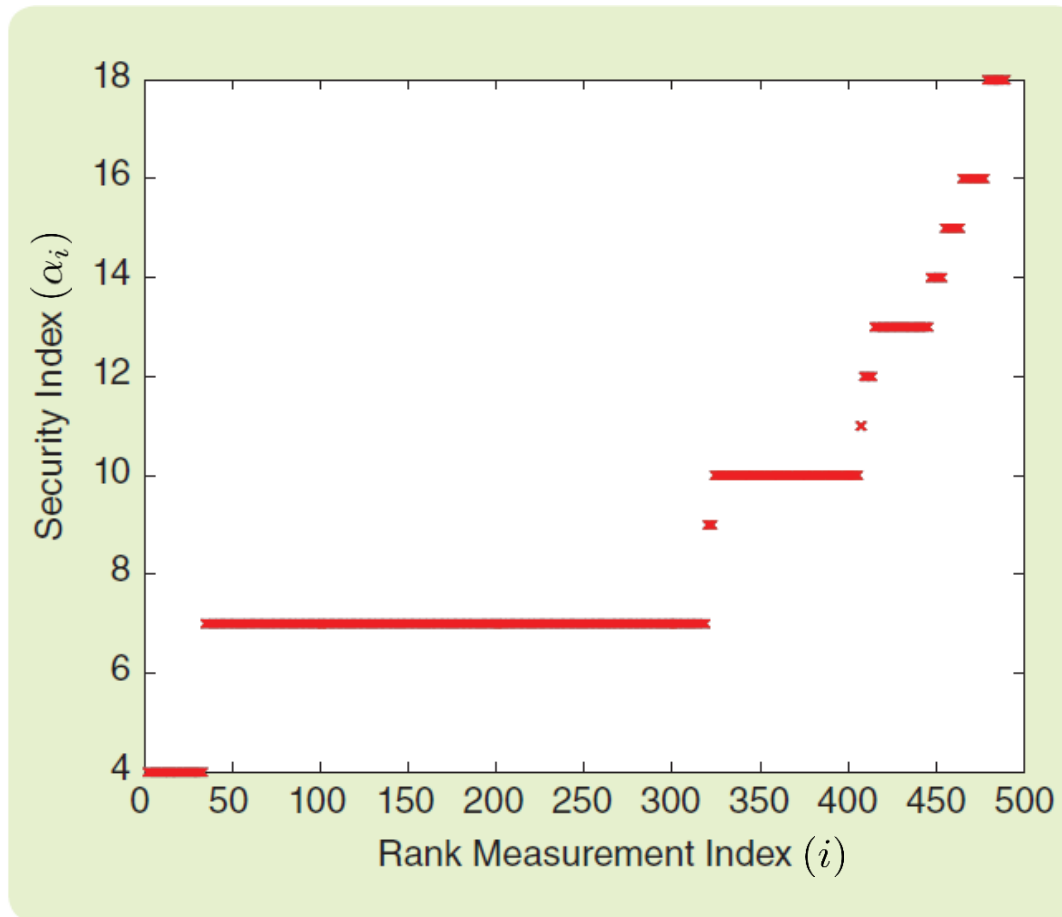
" $\infty$ "

# Example: Power System State Estimator for IEEE 118-bus System

- State dimension  $n = 118$
- Number sensors  $p \approx 490$



# Example: Power System State Estimator for IEEE 118-bus System



- Computation time on laptop using min-cut method [Hendrickx *et al.*, IEEE TAC, 2014]: 0.17 sec
- Used for protection allocation in [Vukovic *et al.*, IEEE JSAC, 2012]



# Summary So Far

- Dynamical security index  $\alpha_i$  defined
- Argued  $\alpha_i$  useful in risk management for assessing likelihood of malicious attack against element  $a_i$
- Computation is NP-hard in general, but often “simple” in special cases:
  - One-dimensional zero-dynamics
  - Static systems with special matrix structures (derived from potential flow problems)
  - Dynamics generally simplifies computation and redundant sensors increase  $\alpha_i$
- Fast computation enables greedy security allocation



# Outline

- Background and motivation
- CPS attack models
- Risk management
- Attack detectability and Security metrics
- **Attack identification and secure state estimation**

# Attack Identification

$$x(k+1) = Ax(k) + B_d d(k) + B_a a(k)$$

$$y(k) = Cx(k) + D_d d(k) + D_a a(k)$$

- Unknown state  $x(k) \in \mathbb{R}^n$
  - Unknown (natural) disturbance  $d(k) \in \mathbb{R}^o$
  - Unknown (malicious) attack  $a(k) \in \mathbb{R}^m$
  - Known measurement  $y(k) \in \mathbb{R}^p$
  - Known model  $A, B_d, B_a, C, D_d, D_a$
- When can we decide there is an attack signal  $a_i \neq 0$ ?
  - Which elements  $a_i$  can we track (“identify”)?
- Not equivalent to designing an unknown input observer/secure state estimator (state not requested here). See end of presentation





# Attack Identification

**Definition:** A (persistent) attack signal  $a$  is

- *identifiable* if for all attack signals  $\tilde{a} \neq a$ , and all corresponding disturbances  $d$  and  $\tilde{d}$ , and initial states  $x(0)$  and  $\tilde{x}(0)$ , we have  $\tilde{y} \neq y$ ;
- *i-identifiable* if for all attack signals  $a$  and  $\tilde{a}$  with  $\tilde{a}_i \neq a_i$ , and all corresponding disturbances  $d$  and  $\tilde{d}$ , and initial states  $x(0)$  and  $\tilde{x}(0)$ , we have  $\tilde{y} \neq y$

## Interpretations:

- Identifiability  $\Leftrightarrow$  (different attack  $a \Rightarrow$  different measurement  $y$ )  $\Leftrightarrow$  attack signal is injectively mapped to  $y \Rightarrow$  attack signal is detectable
- *i-identifiable* weaker than *identifiable*
- $\forall i: a$  is *i-identifiable*  $\Leftrightarrow a$  is *identifiable*
- $a$  is *i-identifiable*: Possible to track element  $a_i$ , but not necessarily  $a_j$ ,  $j \neq i$



# Theorem

Suppose that the attacker can manipulate at most  $q$  attack elements simultaneously ( $\|a\|_0 \leq q$ ).

- i. There exists persistent undetectable attacks  $a^i \Leftrightarrow q \geq \alpha_i$ ;
- ii. All persistent attacks are  $i$ -identifiable  $\Leftrightarrow q < \alpha_i/2$ ;
- iii. All persistent attacks are identifiable  $\Leftrightarrow q < \min_i \alpha_i/2$ .

**Proof.** Compressed sensing type argument. See [Sandberg and Teixeira, SoSCYPS, 2016] for details

## Simple Example of Security Index (cont'd)

$$P(z) = \begin{bmatrix} A - zI & B_d & B_a \\ 0 & 0 & D_a \end{bmatrix} \quad D_a = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Security indices:  $\alpha_1 = 3$ ,  $\alpha_2 = 3$ ,  $\alpha_3 = 3$ ,  $\alpha_4 = \infty$

Attacker with  $q = 1$ : Defender can identify (and thus detect) all attacks

$q = 2$ : Defender can detect (not identify) all attacks against  $a_1, a_2, a_3$  and identify all attacks against  $a_4$

$q = 3 - 4$ : Defender can identify all attacks against  $a_4$ .  
Exist undetectable attacks against  $a_1, a_2, a_3$

## Back to Risk Management

$$P(z) = \begin{bmatrix} A - zI & B_d & B_a \\ 0 & 0 & D_a \end{bmatrix} \quad D_a = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Security indices:  $\alpha_1 = 3$ ,  $\alpha_2 = 3$ ,  $\alpha_3 = 3$ ,  $\alpha_4 = \infty$

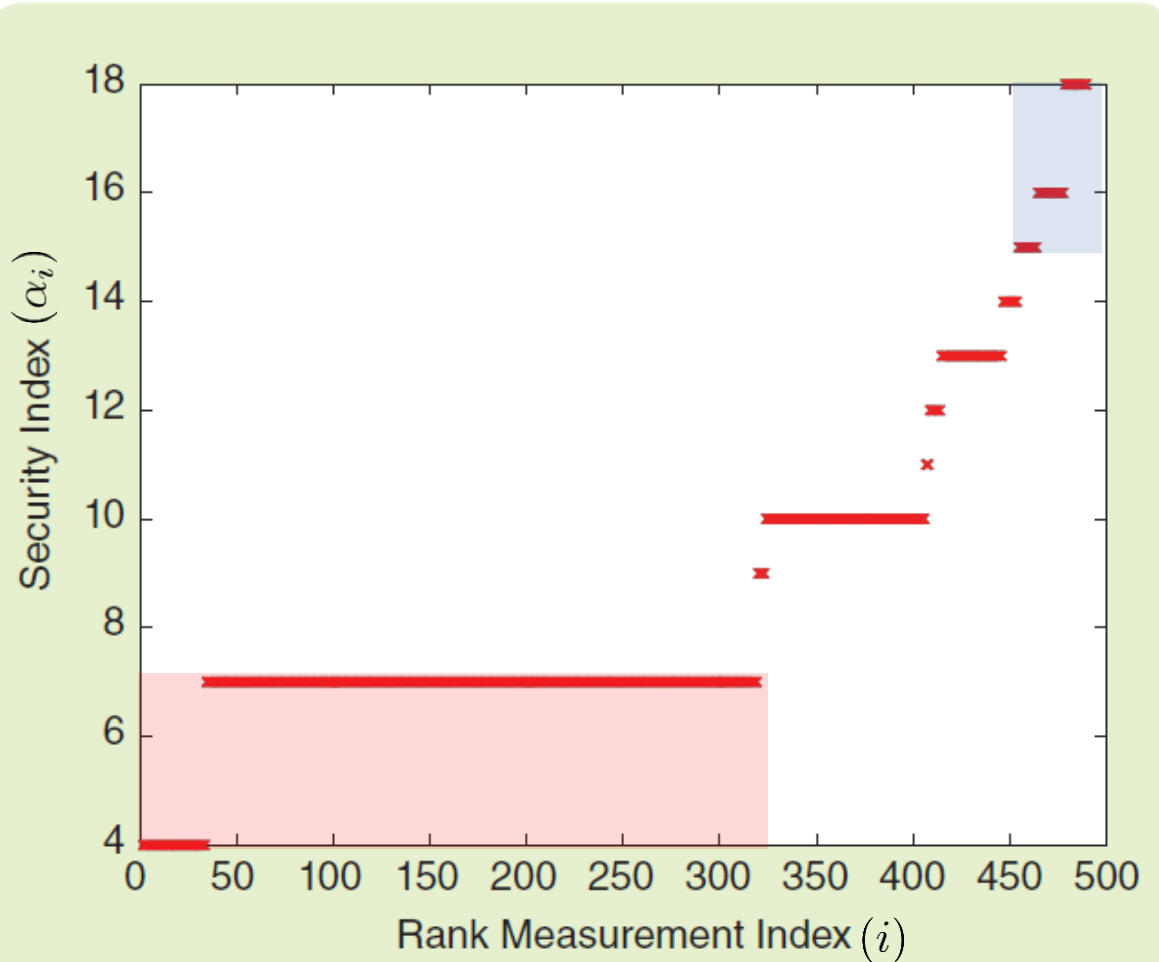
- Suppose the operator can choose to block *one* attack signal (through installing physical protection, authentication, etc.).
- Which signal  $a_1, a_2, a_3$ , or  $a_4$  should she/he choose?
- Among the one(s) with lowest security index! Choose  $a_1$ .
- New attack model and security indices:  $\alpha_2 = \alpha_3 = \alpha_4 = \infty$

$$D_a = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

- By explicitly blocking one attack signal, all other attacks are implicitly blocked (they are identifiable)

# Example: Power System State Estimator for IEEE 118-bus System

- Suppose number of attacked elements is  $q \leq 7$



- Signals susceptible to undetectable attacks

- Signals where all attacks are identifiable

- Other signals will, if attacked, always result in non-zero output  $y$



# Secure State Estimation/Unknown Input Observer (UIO)

**Secure state estimate  $\hat{x}$** : Regardless of disturbance  $d$  and attack  $a$ , the estimate satisfies  $\hat{x} \rightarrow x$  as  $k \rightarrow \infty$

1. Rename and transform attacks and disturbances:

$$\begin{bmatrix} B_d \\ D_d \end{bmatrix} d + \begin{bmatrix} B_a \\ D_a \end{bmatrix} a = \begin{bmatrix} B_f \\ D_f \end{bmatrix} f, \quad \text{such that } \begin{bmatrix} B_f \\ D_f \end{bmatrix} \text{ full column rank}$$

2. Compute security indices  $\alpha_i$  with respect to  $f$

**Theorem:** A secure state estimator exists iff

1.  $(C, A)$  is detectable; and
2.  $q < \min_i \frac{\alpha_i}{2}$ , where  $q$  is max number of non-zero elements in  $f$ .

**Proof.** Existence of UIO by [Sundaram *et al.*, 2007] plus previous theorem

# How to Identify an Attack Signal?

Use decoupling theory from fault diagnosis literature [Ding, 2008]

Suppose that  $y = G_d d + G_a a$  and

$$\begin{aligned}\text{normalrank} [G_d(z)] &= m', \\ \text{normalrank} [G_d(z) \ G_a(z)] &= m' + m''\end{aligned}$$

Then there exists linear decoupling filter  $R$  such that

$$\begin{aligned}\begin{bmatrix} r \\ y' \end{bmatrix} &= R(G_d d + G_a a) = \begin{bmatrix} 0 & \Delta \\ G'_d & G'_a \end{bmatrix} \begin{bmatrix} d \\ a \end{bmatrix}, \\ \text{normalrank} [G'_d(z)] &= \text{normalrank} [G'_d(z) \ G'_a(z)] = m' \\ \text{normalrank} [\Delta(z)] &= m''\end{aligned}$$



# How to Identify an Attack Signal?

Suppose  $a$  is identifiable ( $q < \min_i \alpha_i/2$ )

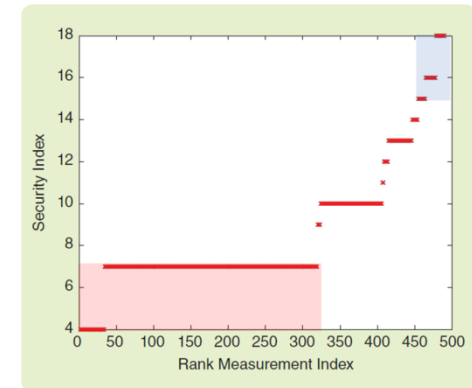
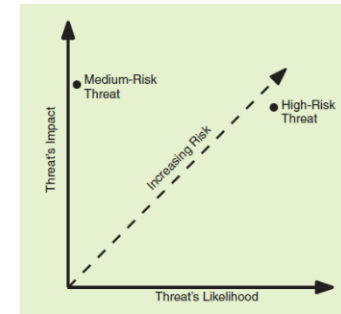
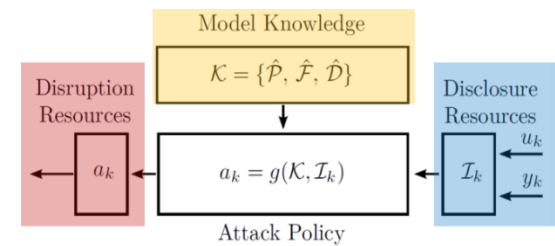
1. Decouple the disturbances to obtain system  $r = \Delta a$
2. Filter out uncertain initial state component in  $r$  to obtain  $r' = \Delta a$
3. Compute left inverses of  $\Delta_I := [\Delta_i]_{i \in I}$  formed out of the columns  $\Delta_i$  of  $\Delta$ , for all subsets  $|I| = q$ ,  $I \subseteq \{1, \dots, m\}$  (**Bottleneck! Compare with compressed sensing**)
4. By identifiability, if estimate  $\hat{a}_I$  satisfies  $r' = \Delta \hat{a}_I$ , then  $\hat{a}_I \equiv a$

(Similar scheme applies if  $a$  is only  $i$ -identifiable)



# Summary

- There is a need for CPS security
- Briefly introduced CPS attack models and concept of risk management
- Input observability and detectability  
⇒ Undetectable attacks and masking initial states and disturbances
- A security metric  $\alpha_i$  for risk management
  - Suppose attacker has access to  $q$  resources:
    - Undetectable attacks against  $a_i$  iff  $q \geq \alpha_i$
    - Attack against  $a_i$  identifiable iff  $q < \alpha_i/2$
- Many useful results in the fault diagnosis literature, especially for identifiable attacks: Unknown input observers, decoupling filters, etc.
- Future research direction: More realistic attacker models, estimate attack likelihoods and impacts, corporation with IT security,...





# Further Reading

## Introduction to CPS/NCS security

- Cardenas, S. Amin, and S. Sastry: "Research challenges for the security of control systems". Proceedings of the 3rd Conference on Hot topics in security, 2008, p. 6.
- Special Issue on CPS Security, IEEE Control Systems Magazine, February 2015
- D. Urbina *et al.*: "Survey and New Directions for Physics-Based Attack Detection in Control Systems", NIST Report 16-010, November, 2016

## CPS attack models, impact, and risk management

- A. Teixeira, I. Shames, H. Sandberg, K. H. Johansson: "A Secure Control Framework for Resource-Limited Adversaries". Automatica, 51, pp. 135-148, January 2015.
- A. Teixeira, K. C. Sou, H. Sandberg, K. H. Johansson: "Secure Control Systems: A Quantitative Risk Management Approach". IEEE Control Systems Magazine, 35:1, pp. 24-45, February 2015
- D. Urbina *et al.*: "Limiting The Impact of Stealthy Attacks on Industrial Control Systems", 23rd ACM Conference on Computer and Communications Security, October, 2016



# Further Reading

## Detectability and identifiability of attacks

- S. Sundaram and C.N. Hadjicostis: “Distributed Function Calculation via Linear Iterative Strategies in the Presence of Malicious Agents”. IEEE Transactions on Automatic Control, vol. 56, no. 7, pp. 1495–1508, July 2011.
- F. Pasqualetti, F. Dörfler, F. Bullo: “Attack Detection and Identification in Cyber-Physical Systems”. IEEE Transactions on Automatic Control, 58(11):2715-2729, 2013.
- H. Fawzi, P. Tabuada, and S. Diggavi: “Secure estimation and control for cyber-physical systems under adversarial attacks”. IEEE Transactions on Automatic Control, vol. 59, no. 6, pp. 1454–1467, June 2014.
- Y. Mo, S. Weerakkody, B. Sinopoli: “Physical Authentication of Control Systems”. IEEE Control Systems Magazine, vol. 35, no. 1, pp. 93-109, February 2015.
- R. Smith: “Covert Misappropriation of Networked Control Systems”. IEEE Control Systems Magazine, vol. 35, no. 1, pp. 82-92, February 2015.
- H. Sandberg and A. Teixeira: “From Control System Security Indices to Attack Identifiability”. Science of Security for Cyber-Physical Systems Workshop, CPS Week 2016



# Further Reading

## Security metrics (security index)

- O. Vukovic, K. C. Sou, G. Dan, H. Sandberg: "Network-aware Mitigation of Data Integrity Attacks on Power System State Estimation". IEEE Journal on Selected Areas in Communications (JSAC), 30:6, pp. 1108--1118, 2012.
- J. M. Hendrickx, K. H. Johansson, R. M. Jungers, H. Sandberg, K. C. Sou: "Efficient Computations of a Security Index for False Data Attacks in Power Networks". IEEE Transactions on Automatic Control: Special Issue on Control of CPS, 59:12, pp. 3194-3208, December 2014.
- H. Sandberg and A. Teixeira: "From Control System Security Indices to Attack Identifiability". Science of Security for Cyber-Physical Systems Workshop, CPS Week 2016

# Acknowledgments

**André M.H. Teixeira**

(Delft University of Technology)



**Kin Cheong Sou**

(National Sun Yat-sen University)



**György Dán**

**Karl Henrik Johansson**

**Jezdimir Milošević**

**David Umsonst**

(KTH)

