

# Formal Impact Metrics for Cyber-physical Attacks

Ruggero Lanotte

DISUIT

Univ. of Insubria, Como, Italy  
ruggero.lanotte@uninsubria.it

Massimo Merro

dept. of Computer Science

Univ. of Verona, Verona, Italy  
massimo.merro@univr.it

Andrei Munteanu

dept. of Computer Science

Univ. of Verona, Verona, Italy  
andrei.munteanu@univr.it

Simone Tini

DISUIT

Univ. of Insubria, Como, Italy  
simone.tini@uninsubria.it

**Abstract**—Cyber-Physical systems (CPSs) are exposed to cyber-physical attacks, i.e., security breaches in cyberspace that adversely affect the physical processes of the systems.

We define two *probabilistic metrics* to estimate the *physical impact* of attacks targeting cyber-physical systems formalised in terms of a probabilistic hybrid extension of Hennessy and Regan’s *Timed Process Language*. Our *impact metrics* estimate the impact of cyber-physical attacks taking into account: (i) the *severity of the inflicted damage* in a given amount of time, and (ii) the *probability* that these attacks are actually accomplished, according to the dynamics of the system under attack. In doing so, we pay special attention to *stealthy attacks*, i.e., attacks that cannot be detected by *intrusion detection systems*. As further contribution, we show that, under precise conditions, our metrics allow us to estimate the impact of attacks targeting a complex CPS in a *compositional way*, i.e., in terms of the impact on its sub-systems.

**Index Terms**—Cyber-physical attacks, impact metrics, timed and hybrid models

## I. INTRODUCTION

*Cyber-Physical Systems* (CPSs) are physical and engineered systems whose operations are monitored, coordinated, controlled, and integrated by a computing and communication core [1]. The growing connectivity and integration of these systems has triggered a dramatic proliferation of *cyber-physical attacks* [2], i.e., security breaches in cyberspace to bring the physical plant into a state desired by the attacker. Some notorious examples are: (i) the *STUXnet* worm, which reprogrammed PLCs of nuclear centrifuges in Iran [3]; (ii) the *CRASHOVER-RIDE* attack on the Ukrainian power grid, otherwise known as *Industroyer* [4]; (iii) the recent *TRITON/TRISIS* malware that targeted a petrochemical plant in Saudi Arabia [5].

Cyber-physical attacks tamper with both the cyber and the physical layer, as they may manipulate both the sensor measurements and the controller commands:

- *attacks on sensors* consist of reading and possibly replacing genuine sensor measurements with fake ones;
- *attacks on actuators* consist of reading, dropping and replacing controller commands with malicious ones.

According to [6], the possible *goals* of cyber-physical attacks can be classified into three main groups: (i) *equipment damage*, such as overstress of equipments, to reduce their expected life cycle, or violation of safety limits; (ii) *production damage*, in order to compromise product quality, product rate or operating costs; (iii) *compliance violations*, such as increasing of environmental pollution.

A clear understanding of both *attacker’s goals* and the *severity of the damages* inflicted when such goals are achieved, is fundamental to conduct a *risk assessment* which ranks

vulnerabilities that may be used by the attacker to achieve its goals. Motivated by the risk assessment application, the objective of this paper is to define formal *probabilistic metrics* to estimate the impact of cyber-physical attacks taking into account both (i) the *severity of the inflicted damage* in a given amount of time, and (ii) the *probability* that these attacks are actually accomplished, according to the dynamics of the CPS under attack. In doing so, we pay special attention to *stealthy attacks*, i.e., attacks that cannot be detected by *intrusion detection systems* (IDSs) which typically monitor sensor signals, control commands, and network communications. In fact, to remain stealthy, attacks would have to closely follow the physical behaviour of the system, causing damage that can be negligible in the short term, but devastating in the long run.

**Contribution:** First of all, we introduce a formal language to specify both CPSs and cyber-physical attacks. For this very purpose, we resort to *process calculi*, a successful and widespread formal approach in *concurrency theory* for representing complex systems such as concurrent, distributed and mobile systems [7], [8], and used in many areas, including verification of security protocols [9] and security analysis of cyber-physical attacks [10].

In Section II, we propose a *hybrid probabilistic process calculus*, inspired by the calculi appeared in [10], [11], and with a clearly-defined *probabilistic behavioural semantics* [12]. In our calculus, cyber-physical systems are represented by making a clear distinction between the *physical component* (also called *plant*) describing the physical process, and the *logical component* describing the cyber devices (i.e., controllers, IDS, supervisors, etc.) governing the physical process as well as the interaction with other cyber devices. At the logical level, our calculus allows us to specify also (MITM) malicious activities targeting physical devices. In particular, we can model integrity attacks on both sensor measurements and actuator commands, as well as drops of actuator commands.

In this formal setting, we design two probabilistic metrics to estimate the impact of (possibly stealthy) cyber-physical attacks. As in Urbina et al. [13], we assume attacks with a complete knowledge of the system under attack, i.e., they know the physical model, the goals in terms of damage that can be inflicted, and the thresholds selected to raise alerts via IDSs.

Our impact metrics are very general, and do not depend on the definition of our calculus, as they rely on the following notions associated to an arbitrary CPS equipped with an IDS:

- a *timed and probabilistic labelled transition semantics* [12] to formally describe the dynamics of both physical and

logical components of (possibly compromised) CPSs;

- a set  $\mathcal{I}$  of *weighted attacker's goal indicators* whose runtime values denote how close is the attacker to reach each of these goals; we assume that IDSs know such indicators (and their associated severities) but they cannot access their values at runtime;
- a *detection policy*  $\mathcal{P}$  that given an *alert signal* (raised by the IDS) and a goal indicator in  $\mathcal{I}$  returns a *statically-determined* estimate of the progresses of the detected attack in achieving the goal denoted by that indicator.

Based on these concepts, in Section III we define two *probabilistic impact metrics*,  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$  and  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$ , to measure the *average false negatives* and the *average false positives*, respectively, faced by a CPS under attack, in the first  $n$  instants of time, according to a weighted attacker's goal indicators  $\mathcal{I}$  and a detection policy  $\mathcal{P}$  associated to the CPS.

Intuitively, given a (possibly compromised) CPS  $M$ , the values  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M) \in [0, 1]$  and  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M) \in [0, 1]$  represent the *average effectiveness* and the *average precision*, respectively, in detecting attacks that reach goals in  $\mathcal{I}$ , raising alert signals (via the IDS of  $M$ ) which are interpreted by the detection policy  $\mathcal{P}$  to estimate the severity of the damages inflicted by those attacks, in the first  $n$  time instants.

Thus, for instance, when  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M) = 0$  the detection policy is very effective in detecting all possible attacks achieving goals in  $\mathcal{I}$ . On the other hand, high values of  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M)$  tell us that the detection policy of  $M$  underestimates the progresses of the attacker in achieving the goals represented via the attacker's goal indicators  $\mathcal{I}$  (high number of *false negatives*).

This metric does not say anything about the precision of the detection: there can be plenty of false positives, thus, apparently, it might seem that a certain attack has inflicted a severe damage, when actually that damage is negligible with respect to  $\mathcal{I}$ .

The average precision of the detection is measured using the metric  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$ . Thus, for instance, when  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M) = 0$  the detection policy  $\mathcal{P}$  is very accurate in signalling true attacks achieving goals in  $\mathcal{I}$ . Whereas, high values of  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M)$  tell us that the detection policy overestimates the progresses of the attacker in achieving the goals represented via the attacker's goal indicators  $\mathcal{I}$  (high number of *false positives*).

*True positive rates* and *false positive rates* have been widely used in the literature to derive metrics, such as ROC curves, for classifying the accuracy of IDSs in terms of the intrusion detection capability, the Bayesian detection rate, and expected cost, which are all multi-criteria optimisation problems between these two rates [14]. In this paper, we take a different point of view. In fact, our metrics will not be used to compare IDSs, but instead for a formal estimation of the impact of (possibly stealthy) cyber-physical attacks, in terms of the physical inflicted damages represented via the attacker's goal indicators  $\mathcal{I}$ , and the logical signals raised by the IDS to detect attacks reaching such goals, according to the detection policy  $\mathcal{P}$ .

As a further contribution, we show that, under precise conditions, our metrics allows us to estimate the impact of cyber-physical attacks on a complex CPS in a *compositional manner*, i.e., in terms of the impact on its sub-systems. In this

respect, our process calculus supports features to compose CPSs in such a way as to avoid interference on physical and/or control variables. Basically, we build up composite CPSs putting together smaller sub-systems whose physical processes remain under the exclusive control of their associated logical components (controllers, IDSs, etc). Thus, interactions between logical components of different sub-systems are always possible, whereas interactions between physical processes may only occur if the corresponding logical components agree on when and how that should happen. Examples of this kind of CPSs can be found in several domains, such as: (i) manufacturing robotic applications (painting, welding, assembly, etc), in which different robotic arms are handled by different, but coordinated, controllers; (ii) drone swarms, where a fleet of drones explores and analyses physical perimeters; (iii) enrichment nuclear facilities, in which series of nuclear centrifuges work under the exclusive control of their programmable logic controllers.

In Section IV, we use our metrics to estimate the impact of four different, carefully chosen, cyber-physical attacks that target a non-trivial use case, consisting of two supervised self-coordinating refrigerated engine systems simulated in MATLAB [15]. Two out of four attacks are actually *stealthy attacks*.

In Section V, we draw conclusions and discuss related work.

## II. A CALCULUS FOR CPSs AND ATTACKS

In this section, we define our *Probabilistic Calculus of Cyber-Physical Systems and Attacks*, called  $\text{pCCPSA}$ , a formal abstract language to specify both physical and logical components of CPSs, as well as cyber-physical attacks. The calculus is inspired by the process calculi appeared in [10], [11]. Unlike those calculi, for the sake of simplicity, we abstract away from physical devices (sensors and actuators) by allowing the cyber components to directly access *observable physical variables* and *control variables*. We also represent *unobservable physical variables*, that cannot be accessed by cyber components.

Let us start with some preliminary notations.

**Notation 1.** We use  $v, w \in \mathcal{V}$  for variables, *partitioned* in observable physical variables  $y, y_k \in \mathcal{Y}$ , unobservable physical variables  $x, x_k \in \mathcal{X}$ , and control variables  $a, a_k \in \mathcal{A}$ . Moreover we use  $c, d \in \mathcal{C}$  for communication channels. Values, ranged over by  $m, m' \in \mathcal{M}$ , are built from basic values, such as integers and real numbers. Given a set of variables  $\mathcal{V}$ , we write  $\mathbb{R}^{\mathcal{V}}$  to denote the set of functions assigning a real value to each variable in  $\mathcal{V}$ . For  $\xi \in \mathbb{R}^{\mathcal{V}}$ ,  $v \in \mathcal{V}$  and  $r \in \mathbb{R}$ , we write  $\xi[v \mapsto r]$  to denote the function  $\psi \in \mathbb{R}^{\mathcal{V}}$  such that  $\psi(w) = \xi(w)$ , for any  $w \neq v$ , and  $\psi(v) = r$ . Given  $\xi_1 \in \mathbb{R}^{\mathcal{V}_1}$  and  $\xi_2 \in \mathbb{R}^{\mathcal{V}_2}$  such that  $\mathcal{V}_1 \cap \mathcal{V}_2 = \emptyset$ , we write  $\xi_1 \uplus \xi_2$  for the function in  $\mathbb{R}^{\mathcal{V}_1 \cup \mathcal{V}_2}$  such that  $(\xi_1 \uplus \xi_2)(v) = \xi_1(v)$ , if  $v \in \mathcal{V}_1$ , and  $(\xi_1 \uplus \xi_2)(v) = \xi_2(v)$ , if  $v \in \mathcal{V}_2$ .

As  $\text{pCCPSA}$  is a probabilistic process calculus, we provide the necessary mathematical machinery for its formal definition.

**Definition 1.** A discrete probability distribution over a set of objects  $\mathcal{O}$  is a function  $\gamma: \mathcal{O} \rightarrow [0, 1]$  with  $\sum_{o \in \mathcal{O}} \gamma(o) = 1$ . The support of  $\gamma$  is the set  $\text{supp}(\gamma) = \{o \in \mathcal{O} : \gamma(o) > 0\}$ .

With  $\mathcal{D}(\mathcal{O})$  we denote the set of all finite-support probability distributions over  $\mathcal{O}$ . For any  $o \in \mathcal{O}$ , the point distribution at  $o$ , written  $\bar{o}$ , assigns probability 1 to  $o$  and 0 to all other objects.

In our language, cyber-physical systems consist of two components: (i) the *physical component*, also called *physics*, describing the physical process, and (ii) the *logical component*, also called *logics*, that interacts with the physics and supports communications between logical sub-components. The physics has two sub-components: the *evolution law*, which describes the time-dependent evolution of the physical process, and the *state of the system*, which is supposed to change at runtime.

**Definition 2 (Physics).** Let  $V$  be a set of variables. The physics of a CPS defined on  $V$  is a pair  $(\varepsilon, \xi_V)$ , where:

- (a)  $\varepsilon$  denotes the evolution law represented as a function of type  $\mathbb{R}^V \rightarrow \mathcal{D}(\mathbb{R}^V)$ ;
- (b)  $\xi_V \in \mathbb{R}^V$  is the state function recording the current state of the variables in  $V$ .

A state function  $\xi_V$  returns the current value associated to each variable in  $V$ . For  $V = X \cup Y \cup A$ , we will abbreviate  $(\xi_V)|_X$  with  $\xi_X$ ,  $(\xi_V)|_Y$  with  $\xi_Y$ , and  $(\xi_V)|_A$  with  $\xi_A$ . An evolution law  $\varepsilon$  models the evolution of the physical system, in which changes on the control variables  $A$  may reflect on both observable variables  $Y$  and unobservable ones  $X$ . Notice that, given a state  $\xi_V$ ,  $\varepsilon(\xi_V)$  returns a *probability distribution over all possible next states* to model the presence of *uncertainty* in our models.

**Example 1 (Physics of a refrigerated engine system).** We provide a running example (inspired from [10]), called *Eng*, in which the temperature of a simple engine is maintained within a specific range by means of a cooling system.

The set of variables  $V = X \cup Y \cup A$  of the system consists of: (i) an observable variable  $temp \in Y$  recording the temperature of the engine and with initial value 95; (ii) six unobservable variables  $past_k \in X$ , for  $1 \leq k \leq 6$ , recording the temperatures in the last 6 time instants, and all initially set to 95; (iii) an unobservable variable  $stress \in X$  denoting the level of stress of the engine, due to its operating conditions, with initial value 0 and maximal value 1; (iv) a control variable  $cool \in A$  to turn on/off the cooling system, with initial value off; (v) a control variable  $speed \in A$  to set the engine speed, ranging over values in the set {slow, half, full}; in normal conditions the engine runs at half power.

Thus, the state function  $\xi_V$  is defined as  $\xi_V = \xi_X \uplus \xi_Y \uplus \xi_A$ , for  $V = X \cup Y \cup A$ , and  $X = \{past_1, \dots, past_6, stress\}$ ,  $Y = \{temp\}$  and  $A = \{cool, speed\}$ ; for the sake of simplicity, we assume  $\xi_A$  to be a mapping where  $\xi_A(cool) \in \{\text{on}, \text{off}\}$  and  $\xi_A(speed) \in \{\text{slow}, \text{half}, \text{full}\}$  such that  $\xi_A(cool) = \text{off}$  when  $\xi_A(cool) \geq 0$ ,  $\xi_A(cool) = \text{on}$  when  $\xi_A(cool) < 0$ ,  $\xi_A(speed) = \text{slow}$  when  $\xi_A(speed) < 0$ ,  $\xi_A(speed) = \text{half}$  when  $\xi_A(speed) = 0$  and  $\xi_A(speed) = \text{full}$  when  $\xi_A(speed) > 0$ .

The evolution law  $\varepsilon$  affects both the observable variables in  $Y$  and the unobservable variables in  $X$  as follows: (i) when the cooling system is active the temperature decreases by a value chosen in the real interval  $[0.8, 1.2]$  according to a discrete probability distribution with granularity  $10^{-1}$ ; (ii) when the

cooling system is inactive the temperature increases by a value  $m$  determined according to a discrete probability distribution with granularity  $10^{-1}$ , and depending on the engine speed: thus, when the speed is slow the value  $m$  is chosen in the real interval  $[0.1, 0.3]$ , when the speed is half then  $m$  is chosen in the real interval  $[0.4, 0.7]$ , when the speed is full then  $m$  is chosen in the real interval  $[0.8, 1.2]$ ; the reader may notice as the uncertainty of the temperature increase grows with the speed of the engine; (iii) the variables  $past_k$ , for  $k \in 1 \dots 6$ , are updated to record the last six temperatures of the system; (iv) the variable  $stress$  remains unchanged if the temperature was in the interval  $[50, 100]$  for at least 3 of the last 6 time instants; otherwise, the variable is increased (reaching at most the value 1) by a constant  $stress\_incr = \frac{1}{50}$ .

Formally, given a state  $\xi_V = \xi_X \uplus \xi_Y \uplus \xi_A$ , the application  $\varepsilon(\xi_V)$  returns a distribution  $\gamma$  such that, for any  $\xi'_X$  and  $\xi'_Y$  with  $\gamma(\xi'_X \uplus \xi'_Y \uplus \xi_A) > 0$ , it holds that:

- If  $\xi_A(cool) = \text{on}$  (active cooling), then  $\gamma(\xi'_X \uplus \xi'_Y \uplus \xi_A) = \frac{1}{5}$  and  $\xi'_Y(temp) = \xi_Y(temp) + m$ , where  $m \in \{-1 + k \cdot 10^{-1} : k \in \mathbb{Z} \wedge -2 \leq k \leq 2\}$ ;
- If  $\xi_A(cool) = \text{off}$  (no cooling) and  $\xi_A(speed) = \text{slow}$ , then  $\gamma(\xi'_X \uplus \xi'_Y \uplus \xi_A) = \frac{1}{3}$  and  $\xi'_Y(temp) = \xi_Y(temp) + m$ , where  $m \in \{0.2 + k \cdot 10^{-1} : k \in \mathbb{Z} \wedge -1 \leq k \leq 1\}$ ;
- If  $\xi_A(cool) = \text{off}$  (no cooling) and  $\xi_A(speed) = \text{half}$ , then  $\gamma(\xi'_X \uplus \xi'_Y \uplus \xi_A) = \frac{1}{4}$  and  $\xi'_Y(temp) = \xi_Y(temp) + m$ , where  $m \in \{0.5 + k \cdot 10^{-1} : k \in \mathbb{Z} \wedge -1 \leq k \leq 2\}$ ;
- If  $\xi_A(cool) = \text{off}$  (no cooling) and  $\xi_A(speed) = \text{full}$ , then  $\gamma(\xi'_X \uplus \xi'_Y \uplus \xi_A) = \frac{1}{5}$  and  $\xi'_Y(temp) = \xi_Y(temp) + m$ , where  $m \in \{1 + k \cdot 10^{-1} : k \in \mathbb{Z} \wedge -2 \leq k \leq 2\}$ ;
- $\xi'_X(past_1) = \xi_X(temp)$  and  $\xi'_X(past_k) = \xi_X(past_{k-1})$ , for  $2 \leq k \leq 6$ ;
- $\xi'_X(stress) = \xi_X(stress)$  if  $|\{k : \xi_X(past_k) \in [50, 100] \wedge 1 \leq k \leq 6\}| \geq 3$ , and  $\xi'_X(stress) = \min(1, \xi_X(stress) + stress\_incr)$ , otherwise.

As regards the formalisation of the *logics* of CPSs in pCCPSA, we build on Hennessy and Regan's *Timed Process Language (TPL)* [16]. We extend TPL with two constructs to read and write variables in  $Y \cup A$ . In this manner, we will model both honest and malicious activities on variables.

**Definition 3 (Logics).** Logical processes are defined as follows:

$$\begin{aligned}
P, Q ::= & \text{nil} \mid \text{tick}.P \mid P \parallel Q \mid [\pi.P]Q \\
& \mid [b]\{P\}, \{Q\} \mid P \setminus c \mid X \mid \text{rec } X.P \\
\pi ::= & \text{read } y(z) \mid \text{write } a\langle m \rangle \mid \text{readAll}(\tilde{z}) \\
& \mid \text{read } a(z) \mid \text{write } y\langle m \rangle \mid \text{rcv } c(z) \mid \text{snd } c\langle m \rangle
\end{aligned}$$

We write *nil* for the *terminated process*. The process  $\text{tick}.P$  sleeps for one time unit and then continues as  $P$ .  $P \parallel Q$  denotes the *parallel composition* of concurrent threads  $P$  and  $Q$ . The process  $[\pi.P]Q$  denotes activities under timeout, meaning that either the activity  $\pi$  succeeds and the process evolves into  $P$ , in the current time unit, or the process timeouts and it evolves into  $Q$ , in the next time unit. Prefix  $\text{read } y(z)$ , for  $y \in Y$ , models a *read of the observable physical variable*  $y$ ; whereas the prefix  $\text{write } a\langle m \rangle$ , for  $a \in A$ , represents a *write on the control variable*  $a$ . These prefixes are executed by the

controller. We also assume a prefix  $\text{readAll}(\tilde{z})$ , executed by IDSs, to read at once all observable physical variables and control variables. Prefixes  $\text{rcv } c(z)$  and  $\text{snd } c(m)$  model input and output activities on channel  $c$ , respectively. Furthermore, prefixes of the form  $\text{read } a(z)$  and  $\text{write } y\langle m \rangle$ , with  $a \in A$  and  $y \in Y$ , denote (MITM) malicious activities targeting physical devices. In particular, with  $\text{read } a(z)$  we model a *drop of an actuator command* associated to the control variable  $a$ ; while  $\text{write } y\langle m \rangle$  represents an *integrity attack on the sensor* associated to the physical variable  $y$ .

**Remark 1.** Notice that neither honest nor malicious code may access to unobservable physical variables in  $\mathcal{X}$ .

The process  $[b]\{P\}, \{Q\}$  is the standard conditional, where  $b$  is a decidable guard. The process  $P \setminus c$  is the standard channel restriction operator. Sometimes we write  $P \setminus \{c_1, c_2, \dots, c_n\}$  to mean  $P \setminus c_1 \setminus c_2 \dots \setminus c_n$ . In processes of the form  $\text{tick}.Q$  and  $[\pi.P]Q$ , the occurrence of  $Q$  is said to be *time guarded*. The process  $\text{rec } X.P$  denotes (time-guarded) recursion. As further notation, we write  $T\{^m/v\}$  for the substitution of the variable  $v$  with the value  $m$  in any expression  $T$ . Similarly, in  $T\{^P/X\}$  the process variable  $X$  is replaced with the process  $P$ .

In the rest of the paper we use the following abbreviations.

**Notation 2.** We write  $\pi.P$  for the process  $\text{rec } X. [\pi.P]X$ , where  $X$  does not occur in  $P$ . We write  $\text{snd } c$  and  $\text{rcv } c$ , when channel  $c$  is used for pure synchronisation.

Everything is in place to define our cyber-physical systems.

**Definition 4** (Cyber-physical system). A cyber-physical system with variables  $V = Y \cup X \cup A$  is given by two components:

- (a) a physical component  $\langle \varepsilon, \xi_V \rangle$ ;
- (b) a logical component  $P$  that may access only observable physical variables and control variables in  $Y \cup A$ , and can communicate, via channels, with other logical components.

We write  $\langle \varepsilon; \xi_V \rangle \bowtie P$  to denote the whole CPS, and use  $M$  and  $N$  to range over CPSs. Sometimes, when the evolution law  $\varepsilon$  is clearly identified, we write  $\xi_V \bowtie P$  instead of  $\langle \varepsilon; \xi_V \rangle \bowtie P$ . CPSs of the form  $\xi_V \bowtie P$  are called *evolution-free*. Finally, we write  $\mathcal{M}$  for the set of all possible CPSs in our calculus.

**Example 2** (A refrigerated engine system). Let us complete the formalisation of our running example  $\text{Eng}$  in  $\text{pCCPSA}$ , by defining its logical component. In Example 1, we already said that the temperature of our engine system is maintained by its controller,  $\text{Ctrl}$ , within a specific range by means of a cooling system. Furthermore, the system is equipped with an IDS that is in charge of checking whether the cooling system is active whenever the temperature is above a certain threshold. If this condition is violated then the IDS tries to mitigate the damages.

The whole system is represented in  $\text{pCCPSA}$  as follows:

$$\text{Eng} = \langle \varepsilon; \xi_V \rangle \bowtie (\text{Ctrl} \parallel \text{IDS}) \setminus \{ins\}$$

where the evolution law  $\varepsilon$  and the state function  $\xi_V$  have been already defined in Example 1, while  $ins$  is a private channel for transmitting instructions on the speed of the engine.

Thus, the logical component of the engine consists of two parallel processes:  $\text{Ctrl}$  that models the controller activity, and  $\text{IDS}$  that implements intrusion detection together with a mild form of mitigation.

Intuitively, at each scan cycle the process  $\text{Ctrl}$  first checks the temperature and then waits for instructions to change the regime of the engine. When the sensed temperature is above 100 degrees, the controller activates the coolant; the cooling activity is maintained for 5 consecutive time instants. Otherwise, if the temperature is not above the threshold 100 then the controller turns off the cooling system. In both cases, before re-starting its scan cycle, the controller waits for instructions/requests to change the engine regime, coming from either its local IDS, via channel  $ins$ , or other IDSs, via channel  $req_{in}$ . In particular, if the controller receives instructions from its local IDS to slow down the engine (because of an anomalous overheating) then it will command so by changing speed. Otherwise, if the local IDS does not see any anomaly, while the IDS of another system is requesting the controller to work at full power, then the process  $\text{Ctrl}$  will command so by acting on variable speed.

The process  $\text{IDS}$  checks whether the cooling system is active when the temperature is above 100 degrees. If this safety condition is violated then: (i) it signals at supervisory level the presence of an overheating anomaly at engine ID, via a channel  $warn$ ; (ii) it sends the controller instructions to slow down the engine, via a channel  $ins$ ; (iii) it requests to other engine systems to run at full power, via a channel  $req_{out}$ , to compensate the lack of performance of its own engine. Of course, such a request of compensation addressed to other engine systems may be accepted or not, depending whether the engine receiving the request is in condition to run at full power or not. However, if the IDS realises that the overheating anomaly is resolved then it asks (both local and external) controllers to reset their engines at half power. Formally,  $\text{Ctrl}$  and  $\text{IDS}$  are defined as follows:

$$\begin{aligned} \text{Ctrl} &= \text{rec } X. \text{tick}. \text{read } temp(z). [z > 100] \{ \text{Cooling} \}, \\ &\quad \{ \text{write } cool\langle off \rangle \}. \text{Check} \} \\ \text{Cooling} &= (\text{write } cool\langle on \rangle. \text{tick})^5. \text{Check} \\ \text{Check} &= \text{rcv } ins(z). \text{rcv } req_{in}(r). [z = slow] \{ \text{write } speed\langle z \rangle. X \}, \\ &\quad \{ \text{write } speed\langle r \rangle. X \} \\ \text{IDS} &= \text{rec } Y. \text{readAll}(\tilde{z}). [z_{temp} > 100 \text{ and } z_{cool} = off] \\ &\quad \{ \text{snd } warn\langle ID, hot \rangle. \text{snd } ins\langle slow \rangle. \text{snd } req_{out}\langle full \rangle. Y \}, \\ &\quad \{ \text{snd } warn\langle ID, ok \rangle. \text{snd } ins\langle half \rangle. \text{snd } req_{out}\langle half \rangle. Y \} \end{aligned}$$

In process  $\text{Cooling}$ , with a small abuse of notation, we write  $(\dots)^5. \text{Check}$  to denote that the process between brackets is repeated 5 consecutive times.

We remark that, for the sake of simplicity, our  $\text{IDS}$  is quite basic: for instance, it does not check whether the temperature is too low. However, it is straightforward to replace our  $\text{IDS}$  with a more sophisticated one, containing more informative tests on process variables and/or on control variables.

#### A. Composing cyber-physical systems

In Example 2, we have modelled an engine system that could be a component of a larger CPS. In this section, we describe how we can safely compose CPSs to build up larger systems

avoiding interferences on physical and/or control variables. The basic idea is the following: we can put together different CPSs provided that their physical components remain under the exclusive control of the corresponding logical components (controllers, IDSs, etc). Said in other words, in a composite system, any possible interaction between physical processes must be “filtered” by the associated logics.

Thus, we will say that two CPSs are *physically-disjoint* if their physics have unrelated variables, whereas their logics may share logical channels to support communications among them. This means that in physically-disjoint CPSs logical components may communicate with each other, whereas (indirect) interactions between physical processes may only occur if the associated logical components agree on when and how that should happen.

**Definition 5** (Physically-disjoint CPSs). *Let  $M_j = \langle \varepsilon_j; \xi_{V_j} \rangle \bowtie P_j$ , for  $j \in \{1, 2\}$  be two CPSs. We say that  $M_1$  and  $M_2$  are physically-disjoint if  $V_1 \cap V_2 = \emptyset$ . In this case, we write  $M_1 \uplus M_2$  to denote the CPS defined as  $\langle (\varepsilon_1 \uplus \varepsilon_2); (\xi_{V_1} \uplus \xi_{V_2}) \rangle \bowtie (P_1 \parallel P_2)$ , where  $\varepsilon_1 \uplus \varepsilon_2$  is the evolution law  $\varepsilon$  such that  $\varepsilon(\xi_{V_1} \uplus \xi_{V_2})(\xi'_{V_1} \uplus \xi'_{V_2}) = \varepsilon_1(\xi_{V_1})(\xi'_{V_1}) \cdot \varepsilon_2(\xi_{V_2})(\xi'_{V_2})$ . The generalisation  $M_1 \uplus \dots \uplus M_n$  to  $n$  CPSs is straightforward.<sup>1</sup>*

Another possible interference-free way of composing sub-systems is to allow interactions between a CPS  $M$  and a *pure-logical component*  $Q$ , that is a process acting only on communication channels and not on variables. This is the case of *supervisory components* interacting with lower-level controllers on (possibly private) communication channels.

**Definition 6** (Parallel composition and restriction). *Let  $M = \langle \varepsilon; \xi_V \rangle \bowtie P$ ,  $Q$  a pure-logical process,  $c$  a channel. We write  $M \parallel Q$  for  $\langle \varepsilon; \xi_V \rangle \bowtie (P \parallel Q)$ , and  $M \setminus c$  for  $\langle \varepsilon; \xi_V \rangle \bowtie (P \setminus c)$ .*

Notice that in the composite system  $M \parallel Q$ , the process  $Q$  cannot interfere with the physical evolution of  $M$ , although it can definitely interact with  $M$  via communication channels, eventually affecting its physical behaviour.

Figure 1 provides an example of what kind of CPS architectures can be represented in pCCPSA by composing sub-systems via the three operators  $\uplus$ ,  $\parallel$  and  $\setminus$ .

**Example 3** (A simple supervised self-coordinating engine system). *We provide a simple supervised self-coordinating engine system consisting of two engines, a left engine,  $Eng^L$ , and a right one,  $Eng^R$ , interacting with a supervisory component that checks the correct functioning of the two engines. The composite system is defined as follows:*

$$Sys = ((Eng^L \uplus Eng^R) \parallel Supervisor) \setminus \{warn\}$$

where both the left and the right engines are properly re-labeled in order to: (i) embed the identities (L or R) of the engines, (ii) distinguish the corresponding variables of the two engines, (iii) fix the channels  $req^L$  and  $req^R$  for requesting full power of the engines L and R, respectively. Formally,

<sup>1</sup>We refer the reader to Notation 1 for the definition of  $\xi_{V_1} \uplus \xi_{V_2}$ .

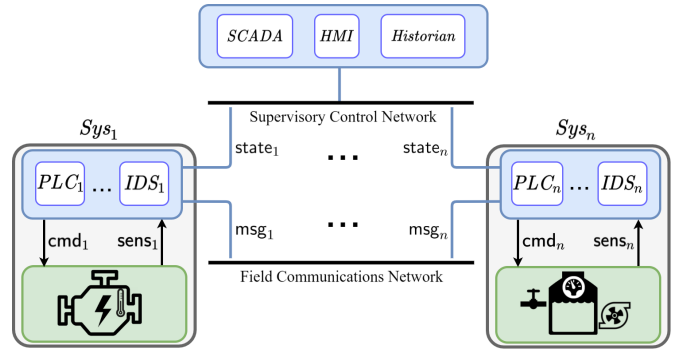


Figure 1. A CPS  $((Sys_1 \uplus \dots \uplus Sys_n) \parallel SCADA \parallel HMI \parallel Hist) \setminus \{state_i\}$

- $Eng^L = Eng \{^L/ID\} \{^{req^L}/req_{in}\} \{^{req^R}/req_{out}\} \{^v/v\}$ , for any  $v \in V$ ;
- $Eng^R = Eng \{^R/ID\} \{^{req^R}/req_{in}\} \{^{req^L}/req_{out}\} \{^v/v\}$ , for any  $v \in V$ .

The pure-logical process Supervisor is defined as follows:

```
rec X.rcv warn(u1, z1).rcv warn(u2, z2).[z1 = z2 = ok]
    {snd okay.tick.X}, {z1 = ok}{snd fail(u2, z2).tick.X},
    {z2 = ok}{snd fail(u1, z1).tick.X}, {snd alarm.tick.X}}
```

Intuitively, if both engines get in a warning state then an alarm signal is sent, otherwise, if only one engine is facing a warning then the supervisor yields a failure to signal which engine is not working properly. Finally, if both engines are working fine then an okay signal is transmitted.

#### B. Probabilistic labelled transition semantics

We start recalling the definition of *probabilistic labelled transition system* (pLTS) [12] over a generic set of terms.

**Definition 7** (pLTS [12]). A probabilistic labelled transition system, or pLTS for short, is a triple  $(\mathcal{T}, \mathbf{A}, \rightarrow)$ , where: (i)  $\mathcal{T}$  is a countable set of terms, (ii)  $\mathbf{A}$  is a countable set of action labels, and (iii)  $\rightarrow \subseteq \mathcal{T} \times \mathbf{A} \times \mathcal{D}(\mathcal{T})$  is a transition relation.

A classical LTS [17] is a special case of a pLTS where all transitions take to point distributions over terms.

Now, we provide the dynamics of pCCPSA in two steps: First, we give a pretty standard LTS for logical processes, then we build up a pLTS for CPSs by lifting transition rules from processes to CPSs, to deal with the probability distributions occurring in evolution laws.

In Table I, we provide the labelled transition semantics for logical components; here, transition rules have the form  $P \xrightarrow{\lambda} Q$ , where  $\lambda$  ranges over actions in the set  $\{\text{tick}, \tau, \bar{c}m, cm, v?m, v!m, (\tilde{z})\}$  denoting: passage of time, unobservable interaction, channel-based transmission/reception, read/write on observable physical variables and/or control variables, respectively, and supervisory read of all variables in  $Y \cup A$ . Rules (Rcv), (Send) and (Comm) model standard channel communication on some channel  $c$ . Rule (Read), for  $v \in Y$ , models the read of the observable physical variable  $v$ ; whereas (Write), for  $v \in A$ , represents the write on the control variable  $v$ . However, for  $v \in A$ , rule (Read) models a malicious attempt of dropping an actuator command acting on the control variable  $v$ . Similarly, when  $v \in Y$ , rule (Write) expresses an *integrity attack*

Table I  
LTS FOR LOGICAL PROCESSES

$\text{(Rcv)} \frac{-}{[\text{rcv } c(z).P]Q \xrightarrow{cm} P\{m/z\}}$	$\text{(Snd)} \frac{-}{[\text{snd } c\langle m \rangle.P]Q \xrightarrow{\bar{c}m} P}$	$\text{(Comm)} \frac{P \xrightarrow{\bar{c}m} P' \quad Q \xrightarrow{cm} Q'}{P \parallel Q \xrightarrow{\tau} P' \parallel Q'}$
$\text{(Read)} \frac{v \in Y \cup A}{[\text{read } v(z).P]Q \xrightarrow{v?m} P\{m/z\}}$	$\text{(Write)} \frac{v \in Y \cup A}{[\text{write } v\langle m \rangle.P]Q \xrightarrow{v!m} P}$	$\text{(Attack)} \frac{P \xrightarrow{v?m} P' \quad Q \xrightarrow{v!m} Q'}{P \parallel Q \xrightarrow{\tau} P' \parallel Q'}$
$\text{(ReadAll)} \frac{-}{[\text{readAll } (\bar{z}).P]Q \xrightarrow{(\bar{z})} P}$	$\text{(TimeNil)} \frac{-}{\text{nil} \xrightarrow{\text{tick}} \text{nil}} \quad \text{(Sleep)} \frac{-}{\text{tick}.P \xrightarrow{\text{tick}} P}$	$\text{(TimeOut)} \frac{-}{[\pi.P]Q \xrightarrow{\text{tick}} Q}$
$\text{(Then)} \frac{\llbracket b \rrbracket = \text{true} \quad P \xrightarrow{\lambda} P'}{\llbracket b \rrbracket \{P\}, \{Q\} \xrightarrow{\lambda} P'}$	$\text{(Else)} \frac{\llbracket b \rrbracket = \text{false} \quad Q \xrightarrow{\lambda} Q'}{\llbracket b \rrbracket \{P\}, \{Q\} \xrightarrow{\lambda} Q'}$	$\text{(Res)} \frac{P \xrightarrow{\lambda} P' \quad \lambda \notin \{cv, \bar{c}v\}}{P \setminus c \xrightarrow{\lambda} P' \setminus c}$
$\text{(Rec)} \frac{P\{\text{rec } X.P/X\} \xrightarrow{\lambda} P'}{\text{rec } X.P \xrightarrow{\lambda} P'}$	$\text{(Par)} \frac{P \xrightarrow{\lambda} P' \quad \lambda \neq \text{tick}}{P \parallel Q \xrightarrow{\lambda} P' \parallel Q}$	$\text{(TimePar)} \frac{P_1 \xrightarrow{\text{tick}} P'_1 \quad P_2 \xrightarrow{\text{tick}} P'_2}{P_1 \parallel P_2 \xrightarrow{\text{tick}} P'_1 \parallel P'_2}$

on (the sensors which reads) the physical variable  $v$ . Thus, in rule (Attack), if  $P$  is a malicious agent then we model the successful dropping of an actuator command emitted by  $Q$ ; whereas if the malicious agent is  $Q$  then the same rule denotes the transmission of a fake value to the agent  $P$  reading the observable physical variable  $v$ . Rule (ReadAll) denotes the reading of all observable physical variables and control variables. Rules (TimeNil), (Sleep) and (TimeOut) model the passage of (discrete) time. Rules (Then), (Else), (Res), (Rec) are standard. Rule (Par) propagates untimed actions over parallel components, whereas rule (TimePar) does the same for timed actions. We omit the symmetric counterparts of (Comm), (Attack) and (Par).

In Table II, we lift the transition rules from logical processes to evolution-free CPSs, by relying on the following notations.

**Notation 3.** We adopt the following notation for probability distributions: given a distribution  $\sigma$  over states and a process  $P$ , we write  $\sigma \bowtie P$  to denote the distribution over (evolution-free) CPSs defined as  $(\sigma \bowtie P)(\xi_V \bowtie P') = \sigma(\xi_V)$ , if  $P = P'$ , and  $(\sigma \bowtie P)(\xi_V \bowtie P') = 0$ , otherwise. Moreover, given an evolution law  $\varepsilon$ , we write  $\langle \varepsilon; \sigma \rangle \bowtie P$  to extend  $\sigma \bowtie P$  to full CPSs as follows:  $(\langle \varepsilon; \sigma \rangle \bowtie P)(\langle \varepsilon'; \xi_V \rangle \bowtie P') = (\sigma \bowtie P)(\xi_V \bowtie P')$ , if  $\varepsilon = \varepsilon'$ , and  $(\langle \varepsilon; \sigma \rangle \bowtie P)(\langle \varepsilon'; \xi_V \rangle \bowtie P') = 0$ , otherwise.

Transition rules have the form  $M \xrightarrow{\alpha} \gamma$ , where  $M$  is an evolution-free CPS and  $\gamma$  is a distribution over (evolution-free) CPSs. For simplicity, as evolution laws contain only static information, the resulting transition rules are parameterised on an evolution law. Thus, Table II provides the transitions rules for evolution-free CPSs.

In Table II, actions, ranged over by  $\alpha$ , are in the set  $\mathbf{A} = \{\tau, \bar{c}m, cm, \text{tick}\}$ . These actions denote: internal activities ( $\tau$ ), channel transmission ( $\bar{c}m$  and  $cm$ ), and the passage of time (tick). Rules (Receive) and (Send) model reception and transmission on a channel  $c$ , respectively. Rule (Tau) lifts silent actions from processes to systems. This includes channel communications and attacks addressed to both physical and control variables. Rule (Read) models the reading of the observable physical variable  $y$ . Notice that the presence of an

attack capable of supplying a fake value (via an action  $y!m$ ) may prevent the access to the correct value by the controller (see rule (Attack) in Table I). Rule (Write) models the writing of a value  $m$  on an control variable  $a$ ; here, the presence of an attacker capable of performing a drop action  $a?m$  may prevent the access to the variable (see, again, rule (Attack) in Table I). Rule (ReadAll) models a read of the current state of all observable physical variables and control variables, this action requires the passage of time as in the next rule, that we explain in more detail. In Rule (Time) timed actions are lifted from logical processes to CPSs; here  $\varepsilon(\xi_V)$  returns a probability distribution over possible state functions for the next time slot, according to the current state  $\xi_V$  and the evolution law  $\varepsilon$ .

Notice that in our pLTS we defined transitions rules of the form  $\xi_V \bowtie P \xrightarrow{\alpha} \sigma \bowtie P'$ , parametric on some evolution law  $\varepsilon$ . As evolution laws do not change at runtime,  $\xi_V \bowtie P \xrightarrow{\alpha} \sigma \bowtie P'$  entails  $\langle \varepsilon; \xi_V \rangle \bowtie P \xrightarrow{\alpha} \langle \varepsilon; \sigma \rangle \bowtie P'$ , thus providing the pLTS for (full) CPSs.

Now, having defined the labelled transitions that can be performed by a CPS of the form  $\langle \varepsilon; \xi_V \rangle \bowtie P$ , we can easily concatenate these transitions to define the possible execution traces of a CPS  $M \in \mathcal{M}$ . Formally,

**Definition 8** (Traces, derivatives and channel communications). We write  $M \xrightarrow{\alpha} M'$  if there exists a transition  $M \xrightarrow{\alpha} \gamma$  such that  $M' \in \text{supp}(\gamma)$ . Let  $t = \alpha_1 \dots \alpha_n$ , for  $n \geq 0$ , a (possibly empty) sequence of actions  $\alpha_i \in \mathbf{A} \setminus \{\tau\}$ . A trace  $M \xrightarrow{t} M'$  is a (possibly empty) sequence of transitions  $M = M_1 \xrightarrow{(\tau)^*} \alpha_1 \xrightarrow{(\tau)^*} \dots \xrightarrow{(\tau)^*} \alpha_n \xrightarrow{(\tau)^*} M_{n+1} = M'$ . Given a CPS  $M \in \mathcal{M}$  we write  $\text{der}(M)$  to denote the set  $\{M' \in \mathcal{M} : \exists t \text{ s.t. } M \xrightarrow{t} M'\}$  of the derivatives of  $M$ . Moreover, we write  $\text{out}(M)$  to denote the set  $\{\bar{c}m \in \mathbf{A} : \exists M' \in \text{der}(M) \text{ s.t. } M' \xrightarrow{\bar{c}m} \gamma\}$  of the signals coming from the (derivatives of)  $M$ , and  $\text{inp}(M)$  to denote the set  $\{cm \in \mathbf{A} : \exists M' \in \text{der}(M) \text{ s.t. } M' \xrightarrow{cm} \gamma\}$  of the possible channel receptions of the (derivatives of)  $M$ .

**Remark 2** (Absence of zeno behaviours). Since recursion is always timed guarded and all timed transitions of CPSs are

Table II  
PLTS FOR CPSS  $\xi_V \bowtie P$  PARAMETRIC ON THE EVOLUTION LAW  $\varepsilon$

$$\begin{array}{l}
\text{(Receive)} \frac{P \xrightarrow{cm} P'}{\xi_V \bowtie P \xrightarrow{cm} \bar{\xi}_V \bowtie P'} \quad \text{(Send)} \frac{P \xrightarrow{\bar{c}m} P'}{\xi_V \bowtie P \xrightarrow{\bar{c}m} \bar{\xi}_V \bowtie P'} \quad \text{(Tau)} \frac{P \xrightarrow{\tau} P'}{\xi_V \bowtie P \xrightarrow{\tau} \bar{\xi}_V \bowtie P'} \\
\text{(Read)} \frac{V = X \cup Y \cup A \quad y \in Y \quad P \xrightarrow{y?m} P' \quad \xi_V(y) = m}{\xi_V \bowtie P \xrightarrow{\tau} \bar{\xi}_V \bowtie P'} \quad \text{(Write)} \frac{V = X \cup Y \cup A \quad a \in A \quad P \xrightarrow{a!m} P'}{\xi_V \bowtie P \xrightarrow{\tau} \bar{\xi}_V[a \mapsto m] \bowtie P'} \\
\text{(ReadIDS)} \frac{P \xrightarrow{(\bar{z})} P' \quad P \xrightarrow{\text{tick}} P \xrightarrow{y?m} \not\rightarrow P \xrightarrow{a!m} \not\rightarrow \quad \tilde{m} = \xi_V(Y \cup A)}{\xi_V \bowtie P \xrightarrow{\text{tick}} \varepsilon(\xi_V) \bowtie P' \{ \tilde{m} / \bar{z} \}} \quad \text{(Time)} \frac{P \xrightarrow{\text{tick}} P' \quad P \xrightarrow{\alpha} \not\rightarrow \quad \alpha \in \{(\bar{z}), y?m, a!m\}}{\xi_V \bowtie P \xrightarrow{\text{tick}} \varepsilon(\xi_V) \bowtie P'}
\end{array}$$

inherited from timed transitions of their logical components (rules (ReadIDS) and (Time)), it follows that the number of untimed transitions between two timed ones is always bounded.

### III. TWO PROBABILISTIC IMPACT METRICS

As said in the Introduction, the two impact metrics  $\mathbf{FN}_{\mathcal{I}, \mathcal{P}}^n$  and  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n$  that we are going to define will be based on the following concepts: (a) a *timed and probabilistic labelled transition semantics* [12] to formally describe the probabilistic dynamics of (possibly compromised) CPSs; (b) a set  $\mathcal{I}$  of *weighted attacker's goal indicators* whose runtime values denote how close is the attacker to reach each of these goals; (c) a *detection policy*  $\mathcal{P}$  that given an alert signal  $\alpha \in \mathbf{A}$  (raised by the IDS) and an attacker's goal indicator in  $\mathcal{I}$  returns a *statically-determined* estimate on the progresses of the detected attack in achieving that goal.

In Table II, we already provided a pLTS to formalise the semantics of CPSs and attacks. Now, let us formally define the two parameters  $\mathcal{I}$  and  $\mathcal{P}$  for a given CPS  $M \in \mathcal{M}$ .

The set of *weighted attacker's goal indicators* is defined as  $\mathcal{I} = \{(i_1, r_1), \dots, (i_k, r_k)\}$ , where  $i_j$  is a function of type  $\text{der}(M) \rightarrow [0, 1]$  and  $r_j \in [0, 1]$ , for  $1 \leq j \leq k$ . Each pair  $(i, r) \in \mathcal{I}$  focuses on some specific *attacker's goal* to compromise the CPS  $M$  under examination and all its derivatives in  $\text{der}(M)$ . Intuitively, for any  $M' \in \text{der}(M)$ ,  $i(M')$  returns a *dynamic quantification* of how close is the attacker injected in  $M'$  to reach the goal  $i$ , whereas the weight  $r$  provides us with a *static information* about the damage inflicted by an attack achieving the goal represented by the indicator  $i$ . Thus, when  $i(M') = 0$  the physical state of  $M'$  is uncorrupted with respect to the attacker's goal indicator  $i$  (e.g., the runtime values of the *stress* variable of the system of Example 2 are close to 0), while  $i(M') = 1$  denotes full corruption of the physical state with respect to the goal indicator  $i$  (e.g., our *stress* variable reaches values close to 1); intermediate values tell us to which extent the physical state of the system  $M'$  is corrupted. In the following, we will write  $\mathcal{I}^- = \{i_1, \dots, i_k\}$  to denote the set of (pure) attacker's goal indicators.

The *detection policy*  $\mathcal{P}: \text{out}(M) \times \mathcal{I}^- \rightarrow [0, 1]$  statically associates alert signals  $\alpha \in \text{out}(M)$  and goal indicators  $i \in \mathcal{I}^-$  to weights that denote an estimate of the progresses of the detected attack, associated to the signal  $\alpha$ , in achieving the goal with indicator  $i$ . Thus, for example, when  $\mathcal{P}(\alpha, i) = 1$

the presence at runtime of the alert signal  $\alpha$  says to system engineers that the attacker's goal denoted by the indicator  $i$  has been fully achieved by the attacker; whereas, when  $\mathcal{P}(\alpha, i) = 0$  the attack detected via the alert signal  $\alpha$  has not done any progress in the achievement of the goal represented by  $i$ . In case of intermediate values, the presence of the signal  $\alpha$  says to which extent the attacker's goal  $i$  has been achieved.

As an example, let us define the weighted attacker's goal indicators and the detection policy for the CPS of Example 2.

**Example 4** (Attacker's goal indicators and detection policy for *Eng*). The set of weighted attacker's goal indicators for *Eng* is defined as follows:  $\mathcal{I}_E = \{(\text{stress}, 1)\}$ . This means that the attacker's goal indicator *stress* is a critical indicator whose variations denote critical attacks. Here, the only attacker's goal indicator coincides with the unobservable physical variable *stress*.<sup>2</sup> The detection policy  $\mathcal{P}_E: \text{out}(\text{Eng}) \times \mathcal{I}_E^- \rightarrow [0, 1]$  is defined as follows:  $\mathcal{P}_E(\overline{\text{warn}}\langle \text{ID}, \text{hot} \rangle, \text{stress}) = 1$  and  $\mathcal{P}_E(\alpha, i) = 0$ , otherwise. This means that the presence of the signal  $\overline{\text{warn}}\langle \text{ID}, \text{hot} \rangle$  emitted by the IDS provides an estimate that the attacker fully achieved the goal of affecting the stress of the engine system. On the other hand, transmissions at communications channel  $\text{req}_{\text{out}}$  are considered irrelevant from the point of view of the attacker's goal indicators  $\mathcal{I}_E^-$ .

In order to define our impact metrics  $\mathbf{FN}_{\mathcal{I}, \mathcal{P}}^n$  and  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n$ , we define a new probabilistic transition  $M \xrightarrow{q_1, \dots, q_k} \mathcal{P} \gamma$  summing up "attacker's progresses" associated to execution traces spanning over a single time slot, according to the detection policy  $\mathcal{P}$ . Intuitively, the weights  $q_j \in [0, 1]$ , for  $j \in \{1, \dots, k\}$ , sums up estimated progresses achieved by the attacker in one time slot with respect to each attacker's goal indicators  $i_j$ , according to the alert signals emitted by the IDS. Formally,

**Definition 9.** Let  $M \in \mathcal{M}$  be a CPS with an associated set of weighted attacker's goal indicators  $\mathcal{I}$ , such that  $\mathcal{I}^- = \{i_1, \dots, i_k\}$ , and detection policy  $\mathcal{P}$ . We write  $M \xrightarrow{q_1, \dots, q_k} \mathcal{P} \gamma$  if  $M \xrightarrow{t} M'$ , for a sequence of untimed actions  $t = \alpha_1 \dots \alpha_n$ ,<sup>3</sup> with  $\alpha_1, \dots, \alpha_n \in \mathbf{A} \setminus \{\tau, \text{tick}\}$ , such that  $M' \xrightarrow{\text{tick}} \gamma$  and  $q_j = \min(1, \sum_{l=1, \alpha_l \in \text{out}(M)} \mathcal{P}(\alpha_l, i_j))$ , for  $j \in \{1, \dots, k\}$ .<sup>4</sup>

<sup>2</sup>We recall that physical variables can be seen as functions that given the current state of the CPS return the current value of the variable.

<sup>3</sup>We recall that untimed actions always lead to point distributions.

<sup>4</sup>By Remark 2 the number of untimed actions preceding a tick is always finite.

Table III  
DEFINITIONS OF THE METRICS  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M)$  AND  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M)$  FOR  $n > 0$ .

$$\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M) = \min \left( 1, \max_{M \xrightarrow{q_1, \dots, q_k} \mathcal{P} \gamma} \frac{\sum_{j=1}^k r_j \cdot \max(0, i_j(M) - q_j) + (n-1) \cdot \sum_{M' \in \text{supp}(\gamma)} \gamma(M') \cdot \mathbf{FN}_{\mathcal{I},\mathcal{P}}^{n-1}(M')}{n} \right)$$

$$\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M) = \min \left( 1, \max_{M \xrightarrow{q_1, \dots, q_k} \mathcal{P} \gamma} \frac{\sum_{j=1}^k r_j \cdot \max(0, q_j - i_j(M)) + (n-1) \cdot \sum_{M' \in \text{supp}(\gamma)} \gamma(M') \cdot \mathbf{FP}_{\mathcal{I},\mathcal{P}}^{n-1}(M')}{n} \right)$$

Now, we are ready to define the two metrics  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n: \mathcal{M} \rightarrow [0, 1]$  and  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n: \mathcal{M} \rightarrow [0, 1]$  to measure the *average effectiveness* and the *average precision*, respectively, of the detection mechanism of a (possibly compromised) CPS  $M \in \mathcal{M}$  in recognising attacks achieving goals in  $\mathcal{I}$ , by raising alert signals (via the IDS) that are interpreted by the detection policy  $\mathcal{P}$  to estimate the severity of the damages inflicted by those attacks, in the first  $n$  time instants.

Basically, for each (probabilistic) execution trace leading from a system  $M$  to a system  $M'$  in at most  $n$  time instants, the metric  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$  measures the *effectiveness* of the detection by computing the average (weighted) deviation between the actual damage inflicted by attackers with goals indicators  $i_j \in \mathcal{I}^-$ , and the estimate returned by the policy  $\mathcal{P}$  of those damages. Such a deviation is weighted using both the severity  $r_j$  of the damage associated to  $i_j$  and the probability to reach  $M'$  from  $M$ . Thus, for example, high values of  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M)$  tell us that the detection mechanism underestimates the progresses of the attacker in achieving the attacker's goal indicators represented in  $\mathcal{I}$  (high number of *false negatives*). The definition of  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$  is by induction on the time  $n \in \mathbb{N}$ . Precisely,  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^0(M) = 0$ , while, for  $n > 0$ , the definition is given in Table III and relies on all derivatives  $M'$  of  $M$ , reached in at most  $n$  time instants.

As for the metric  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$ , for each (probabilistic) execution trace leading from a system  $M$  to a system  $M'$  in at most  $n$  time instants,  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$  measures the *precision* of the detection by computing the average (weighted) deviation between the estimate provided by the policy  $\mathcal{P}$  of the damage inflicted by attackers with goals in  $\mathcal{I}$ , and the actual damage achieved by the attackers, according to each weighted attacker's goal indicator  $(i_j, r_j) \in \mathcal{I}$ . Such a deviation is weighted using both the severity  $r_j$  of the damage associated to the indicator  $i_j$  and the probability to reach  $M'$  from  $M$ . Thus, for example, high values of  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M)$  tell us that the detection mechanism overestimates the progresses of the attacker in achieving the attacker's goal indicators represented in  $\mathcal{I}$  (high number of *false positives*). The definition of  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$  is by induction on the time  $n \in \mathbb{N}$ . Precisely,  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^0(M) = 0$ , while, for  $n > 0$ , the definition is given in Table III and relies on all derivatives of  $M$ , reached in at most  $n$  time instants.

As an example, we can apply our metrics to formally state that the detection policy of the engine system  $Eng$  of Example 2 is effective and precise, i.e., there are neither false negatives nor false positives, with respect to the attacker's goal indicators  $\mathcal{I}_E$  and the detection policy  $\mathcal{P}_E$ , as defined in Example 4.

**Proposition 1.**  $\mathbf{FN}_{\mathcal{I}_E, \mathcal{P}_E}^n(Eng) = 0$  and  $\mathbf{FP}_{\mathcal{I}_E, \mathcal{P}_E}^n(Eng) = 0$ .

We conclude this section with a remark on how our metrics deal with nondeterministic behaviours.

**Remark 3** (Worst-case approach). *Although both metrics work probabilistically, when dealing with nondeterministic behaviours of the compromised system (e.g., a possible drop of an actuator command achieving one of the goals in  $\mathcal{I}$ ) the metrics adopt a worst-case approach. More precisely, the metric  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$  chooses the behaviour of the system that most underestimate the attacker's progress; whereas the metric  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$  chooses the behaviour that most overestimate the attacker's progress.*

#### A. Composing the impacts on CPS sub-systems

In Section II-A, we described how to compose larger CPSs putting together physically-disjoint components and pure-logical ones. In this section, we describe how to compute the impact of an attack on a composite CPS in terms of the impacts on its sub-systems. In doing so, we have to understand how to compose the attacker's goal indicators and the detection policies of sub-components. For instance, when putting together two physically-disjoint components  $M_1$  and  $M_2$  we must define how to compose their corresponding attacker's goal indicators  $\mathcal{I}_1$  and  $\mathcal{I}_2$  and the detection policies  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , according to some notion of *security clearance* associated to the each component. Intuitively, the security clearance of a physically-disjoint CPS component should denote the importance of that component from a security point of view (e.g., in an airplane, securing the engine control system may be more critical than securing the air conditioning system).

**Definition 10** (Distributing  $\mathcal{I}$  and  $\mathcal{P}$  over  $\uplus$ ). *Let  $M_1$  and  $M_2$  be two physically-disjoint CPSs, with disjoint weighted attacker's goal indicators  $\mathcal{I}_1$  and  $\mathcal{I}_2$ , and detection policies  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , respectively. Let  $p_1, p_2 \in [0, 1]$  be the security clearances associated to  $M_1$  and  $M_2$ , respectively. Then, the CPS  $M_1 \uplus M_2$  has attacker's goal indicators  $p_1 \cdot \mathcal{I}_1 \uplus p_2 \cdot \mathcal{I}_2$ , and detection policy  $p_1 \cdot \mathcal{P}_1 \uplus p_2 \cdot \mathcal{P}_2$ , where:*

- $p_1 \cdot \mathcal{I}_1 \uplus p_2 \cdot \mathcal{I}_2 = \bigcup_{(i,r) \in \mathcal{I}_1} (i, p_1 \cdot r) \cup \bigcup_{(i,r) \in \mathcal{I}_2} (i, p_2 \cdot r)$ ;
- $p_1 \cdot \mathcal{P}_1 \uplus p_2 \cdot \mathcal{P}_2$  is such that  $(p_1 \cdot \mathcal{P}_1 \uplus p_2 \cdot \mathcal{P}_2)(\alpha, i) = \mathcal{P}_1(\alpha, i)$ , if  $i \in \mathcal{I}_1^-$ , and  $(p_1 \cdot \mathcal{P}_1 \uplus p_2 \cdot \mathcal{P}_2)(\alpha, i) = \mathcal{P}_2(\alpha, i)$ , if  $i \in \mathcal{I}_2^-$ .

*The generalisation to CPSs as  $M_1 \uplus \dots \uplus M_n$  is straightforward.*

A similar definition is necessary for both channel restriction and parallel composition between CPSs and pure-logical processes (we recall that we use pure-logical processes to represent *supervisory components* interacting with lower-level controllers



on communication channels). Note that we have not associated yet attacker's goal indicators and detection policies to pure-logical processes. However, as in a system  $M \parallel Q$  the pure-logical process  $Q$  may communicate with  $M$ , it is reasonable to assume that it may emit signals about the indicators of  $M$ .

**Definition 11** (Distributing  $\mathcal{I}$  and  $\mathcal{P}$  over  $\parallel$  and restriction). *Let  $M \in \mathcal{M}$  be a CPS with attacker's goal indicators  $\mathcal{I}$  and detection policy  $\mathcal{P}$ . Let  $Q$  be a pure-logical process with the following detection policy  $\mathcal{P}_Q: \text{out}(Q) \times \mathcal{I}^- \rightarrow [0, 1]$ , wrt  $\mathcal{I}$ , such that for any  $(\alpha, i) \in \text{dom}(\mathcal{P}) \cap \text{dom}(\mathcal{P}_Q)$  it holds that  $\mathcal{P}(\alpha, i) = \mathcal{P}_Q(\alpha, i)$ . Let  $c$  be a channel. Then,*

- $M \parallel Q$  has weighted attacker's goal indicators  $\mathcal{I}$  and detection policy  $\mathcal{P} \parallel \mathcal{P}_Q$ , where  $(\mathcal{P} \parallel \mathcal{P}_Q)(\alpha, i) = \mathcal{P}(\alpha, i)$ , if  $(\alpha, i) \in \text{dom}(\mathcal{P})$ ;  $(\mathcal{P} \parallel \mathcal{P}_Q)(\alpha, i) = \mathcal{P}_Q(\alpha, i)$ , otherwise, and
- $M \setminus c$  has weighted attacker's goal indicators  $\mathcal{I}$  and detection policy  $(\mathcal{P} \setminus c)$ , where  $(\mathcal{P} \setminus c)(\alpha, i) = \mathcal{P}(\alpha, i)$  if  $c$  does not occur in  $\alpha$ , and  $(\mathcal{P} \setminus c)(\alpha, i) = 0$ , otherwise.

**Example 5** (Attacker's goal indicators and detection policy for the system of Example 3). *Let's define the attacker's goal indicators  $\mathcal{I}_{Sys}$  and the detection policy  $\mathcal{P}_{Sys}$  for the composite system  $Sys$  of Example 3. We derive them from the components  $Eng^L$ ,  $Eng^R$  and Supervisor (for brevity, here, we will call it  $Super$ ). The security clearances of the two engine systems are:  $p_1 = p_2 = 0.7$  (we assume equal criticality). Thus,  $\mathcal{I}_{Sys} = 0.7 \cdot \mathcal{I}_L \uplus 0.7 \cdot \mathcal{I}_R$  and  $\mathcal{P}_{Sys} = ((0.7 \cdot \mathcal{P}_L \uplus 0.7 \cdot \mathcal{P}_R) \parallel \mathcal{P}_{Super}) \setminus \{\text{warn}\}$ , where: (i) the goal indicators  $\mathcal{I}_L$  and  $\mathcal{I}_R$  for  $Eng^L$  and  $Eng^R$ , respectively, are obtained from the goal indicators  $\mathcal{I}_E$  in Example 4, in the obvious manner; (ii) the same holds for the two policies  $\mathcal{P}_L$  and  $\mathcal{P}_R$ ; (iii) the detection policy  $\mathcal{P}_{Super}$  is defined as follows:*

- $\mathcal{P}_{Super}(\overline{\text{alarm}}, \text{stressL}) = \mathcal{P}_{Super}(\overline{\text{alarm}}, \text{stressR}) = 1$ ;
- $\mathcal{P}_{Super}(\overline{\text{fail}}(L, \cdot), \text{stressL}) = \mathcal{P}_{Super}(\overline{\text{fail}}(R, \cdot), \text{stressR}) = 1$ ;
- $\mathcal{P}_{Super}(\overline{\text{fail}}(L, \cdot), \text{stressR}) = \mathcal{P}_{Super}(\overline{\text{fail}}(R, \cdot), \text{stressL}) = 0$ ;
- $\mathcal{P}_{Super}(\overline{\text{okay}}, \text{stressL}) = \mathcal{P}_{Super}(\overline{\text{okay}}, \text{stressR}) = 0$ ;
- $\mathcal{P}_{Super}(\alpha, i) = 0$ , otherwise.

Now, everything is in place to formally state in which terms our impact metrics are compositional with respect to the notion of compositionality given in Definitions 5 and 6 of Section II-A.

We start by computing our metrics on composite systems consisting of physically-disjoint CPSs, under the hypothesis that channel communications among those CPSs do not carry any information on the achievement of attacker's goal indicators, unlike alert signals that are addressed to the supervisory level.

Basically, the numeric values of the impact metrics of a composite system of the form  $M_1 \uplus \dots \uplus M_k$  are obtained through a weighted sum of the impacts on the single components.

**Theorem 1** (Compositionality w.r.t.  $\uplus$ ). *Let  $M_i \in \mathcal{M}$ , for  $i \in \{1, \dots, k\}$ , be  $k$  physically-disjoint CPSs, with security clearances  $p_i \in [0, 1]$ , disjoint attacker's goal indicators  $\mathcal{I}_i$ , and detection policies  $\mathcal{P}_i$ , respectively, such that for any  $i, j \in \{1, \dots, k\}$ ,  $i \neq j$ , whenever  $\bar{c}m \in \text{out}(M_i)$  and  $cm \in \text{inp}(M_j)$  it holds that  $\mathcal{P}_i(\bar{c}m, h) = 0$ , for any  $h \in \mathcal{I}_i^-$ . Then,*

- 1)  $\mathbf{FN}_{\mathcal{I}, \mathcal{P}}^n(M_1 \uplus \dots \uplus M_k) = \min(1, p_1 \cdot \mathbf{FN}_{\mathcal{I}_1, \mathcal{P}_1}^n(M_1) + \dots + p_k \cdot \mathbf{FN}_{\mathcal{I}_k, \mathcal{P}_k}^n(M_k))$ , when  $\mathcal{I} = p_1 \cdot \mathcal{I}_1 \uplus \dots \uplus p_k \cdot \mathcal{I}_k$  and  $\mathcal{P} = p_1 \cdot \mathcal{P}_1 \uplus \dots \uplus p_k \cdot \mathcal{P}_k$  and  $\sum_{(i,r) \in \mathcal{I}_j} r \leq 1$ , for any  $j \in \{1, \dots, k\}$ ;
- 2)  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M_1 \uplus \dots \uplus M_k) = \min(1, p_1 \cdot \mathbf{FP}_{\mathcal{I}_1, \mathcal{P}_1}^n(M_1) + \dots + p_k \cdot \mathbf{FP}_{\mathcal{I}_k, \mathcal{P}_k}^n(M_k))$ , when  $\mathcal{I} = p_1 \cdot \mathcal{I}_1 \uplus \dots \uplus p_k \cdot \mathcal{I}_k$  and  $\mathcal{P} = p_1 \cdot \mathcal{P}_1 \uplus \dots \uplus p_k \cdot \mathcal{P}_k$  and  $\sum_{(i,r) \in \mathcal{I}_j} r \leq 1$ , for any  $j \in \{1, \dots, k\}$ .

When composing CPSs with pure-logical processes (up to eventual channel restrictions) we get, in general, a weaker result. This is because the introduction of a parallel pure-logical process potentially increases the number of (observable) alert signals, whereas channel restriction, in general, reduces their number as it makes signals non-observable.

**Proposition 2** (Compositionality w.r.t.  $\parallel$  and restriction). *Let  $M \in \mathcal{M}$  be a CPS with attacker's goal indicators  $\mathcal{I}$ , and detection policy  $\mathcal{P}$ . Let  $Q$  be a pure-logical process with detection policy  $\mathcal{P}_Q$  over goal indicators  $\mathcal{I}$  such that whenever  $\bar{c}m \in \text{out}(M)$  and  $cm \in \text{inp}(Q)$  then  $\mathcal{P}(\bar{c}m, h) = 0$ , for any  $h \in \mathcal{I}^-$ . Let  $c$  be a channel. Then,*

- 1)  $\mathbf{FN}_{\mathcal{I}, \mathcal{P}'}^n(M \parallel Q) \leq \mathbf{FN}_{\mathcal{I}, \mathcal{P}}^n(M)$ , for  $\mathcal{P}' = (\mathcal{P} \parallel \mathcal{P}_Q)$ ;
- 2)  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}'}^n(M \parallel Q) \geq \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M)$ , for  $\mathcal{P}' = (\mathcal{P} \parallel \mathcal{P}_Q)$ ;
- 3)  $\mathbf{FN}_{\mathcal{I}, \mathcal{P}'}^n(M \setminus c) \geq \mathbf{FN}_{\mathcal{I}, \mathcal{P}}^n(M)$ , for  $\mathcal{P}' = (\mathcal{P} \setminus c)$ ;
- 4)  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}'}^n(M \setminus c) \leq \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M)$ , for  $\mathcal{P}' = (\mathcal{P} \setminus c)$ .

However, in a composite system, when all alert signals  $\mathcal{A} \subseteq \text{out}(M)$  coming from a CPS  $M$  in a single time slot are intercepted by a supervisor (logical-process)  $Sup$  which analyses and aggregates those signals, converting them in alert signals  $\mathcal{B} \subseteq \text{out}(Sup)$  in the same time slot, then we can strengthen the results of Proposition 2. As an example, suppose that  $\mathcal{P}$  and  $\mathcal{P}_{Sup}$  are the detection policy of  $M$  and  $Sup$ , respectively, with respect to the same set  $\mathcal{I}$  of attacker's goal indicators. And suppose that at each time slot the detection policy of the supervisor amplifies the estimate of the attacker's progress made by the detection policy of  $M$ , i.e.,  $\min(1, \sum_{\beta \in \mathcal{B}} \mathcal{P}_{Sup}(\beta, i)) \geq \min(1, \sum_{\alpha \in \mathcal{A}} \mathcal{P}(\alpha, i))$ , for any  $i \in \mathcal{I}$ , abbreviated  $\mathcal{P}_{Sup} \geq \mathcal{P}$ .<sup>5</sup> Then, we can prove that the detection policy of the supervised system  $(M \parallel Sup) \setminus \text{chn}(\mathcal{A})$ <sup>6</sup> is at least as aggressive as the detection policy  $\mathcal{P}$  of the CPS  $M$ . As expected, the increased aggressiveness of the detection policy results in (a possibly) smaller number of false negatives and in (a possibly) greater number of false positives in the composite supervised system. More generally,

**Theorem 2** (Compositionality on supervised systems). *Let  $M \in \mathcal{M}$  be a CPS with attacker's goal indicators  $\mathcal{I}$ , alert signals  $\mathcal{A} \subseteq \text{out}(M)$ , and detection policy  $\mathcal{P}$ . Let  $Sup$  be a pure-logical process, with detection policy  $\mathcal{P}_{Sup}$ , collecting all alert signals in  $\mathcal{A}$  and converting them in alert signals in  $\mathcal{B} \subseteq \text{out}(Sup)$ , on fresh<sup>7</sup> communication channels. If  $\mathcal{P}_{Sup} \sim \mathcal{P}$ , where  $\sim$  is an operator among  $\{\leq, =, \geq\}$ , then:*

<sup>5</sup>The reader is referred to the appendix for a formal definition of detection policy comparison, such as  $\mathcal{P}_{Sup} \geq \mathcal{P}$ ,  $\mathcal{P}_{Sup} = \mathcal{P}$  or  $\mathcal{P}_{Sup} \leq \mathcal{P}$ .

<sup>6</sup>Here,  $\text{chn}(\mathcal{A})$  denotes the set of channels occurring in the actions of  $\mathcal{A}$ .

<sup>7</sup>Here, "fresh" means that the channels used in  $\mathcal{B}$  do not occur elsewhere.

- 1)  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(M) \sim \mathbf{FN}_{\mathcal{I},\mathcal{P}'}^n((M \parallel \text{Sup}) \setminus \text{chn}(\mathcal{A}))$ , for  $\mathcal{P}' = (\mathcal{P} \parallel \mathcal{P}_{\text{Sup}}) \setminus \text{chn}(\mathcal{A})$
- 2)  $\mathbf{FP}_{\mathcal{I},\mathcal{P}'}^n((M \parallel \text{Sup}) \setminus \text{chn}(\mathcal{A})) \sim \mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(M)$ , for  $\mathcal{P}' = (\mathcal{P} \parallel \mathcal{P}_{\text{Sup}}) \setminus \text{chn}(\mathcal{A})$ .

Based on the compositionality results stated in Theorem 1 and Theorem 2, we can lift the result of Proposition 1 to our composite system  $Sys$  of Example 3.

**Proposition 3.** *Let  $\mathcal{I}_{Sys}$  and  $\mathcal{P}_{Sys}$  be as defined in Example 5. Let  $\mathcal{I} = \mathcal{I}_{Sys}$  and  $\mathcal{P} = \mathcal{P}_{Sys}$ . Then,  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(Sys) = 0$  and  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(Sys) = 0$ , for any  $n \in \mathbb{N}$ .*

#### IV. OUR IMPACT METRICS AT WORK

In this section, we analyse four attacks targeting the  $Eng$  system given in Example 2, and aiming at causing either *overstress of the engine* [6] and/or *deception of the IDS* (increasing the false positives) in a possibly stealthy manner.

In order to test the (computational) feasibility of our metrics  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n$  and  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$ , we simulate in MATLAB [15] the  $Eng$  system together with the four attacks mentioned above, computing their impacts during *attack windows* whose duration varies between 300 and 500 time instants. We rely on the compositionality results of Theorem 1 and Theorem 2 to estimate the impact of each attack, for the same attack window, on the composite system  $Sys$  of Example 3.

Our analyses are conducted on a notebook with the following set-up: 2.8 GHz Intel i7 7700 HQ, with 16 GB memory (plus 48 GB for swap), and Ubuntu 20.04 LTS OS. For each attack, the computation of the two metrics required at most 14 hours. The source files of both the  $Eng$  system and the attacks, together with analysis results and scripts to reproduce them, can be found at: [https://bitbucket.org/formal\\_projects/impact\\_metrics/](https://bitbucket.org/formal_projects/impact_metrics/).

The *first attack* is an integrity attack to the actuator of the  $Eng$  system. It consists of forging fake actuator commands to turn off the cooling system. Formally,

$$Att_1 = \text{rec } X.\text{read } temp(y).[y \leq 99.5]\{\text{write } cool\langle\text{off}\rangle.\text{tick}.X\}, \{\text{tick}.X\}$$

Here, the attacker turns off the cooling system as soon as the temperature reaches 99.5 degrees, even before the completion of the 5-ticks cooling cycle. In this manner, the temperature of the engine rises quickly above the threshold 100, accumulating stress in the system, and requiring continuous activations of the cooling system. This is a *stealthy attack* (i.e., undetected by our IDS) as it keeps the system close to the threshold 100 without ever exceeding it; only the attacker's goal indicator  $stress$  keeps track of the physical damage inflicted by the attack.

The impact of the attack in the first 500 time instants is represented in Figure 2. As the reader can see, the average number of false negatives of the  $Eng$  system under attack, i.e.,  $\mathbf{FN}_{\mathcal{I}_E,\mathcal{P}_E}^n(Eng \parallel Att_1)$ , grows with the size of the attack window. Thus, after 500 time instants the metric  $\mathbf{FN}_{\mathcal{I}_E,\mathcal{P}_E}^n$  approaches the value 0.6, meaning that, in average, the attack achieves around 60% of its goals, according to  $\mathcal{I}_E$ , i.e., the goal indicator  $stress$ . Notice that this attack does not affect the false positives at all.

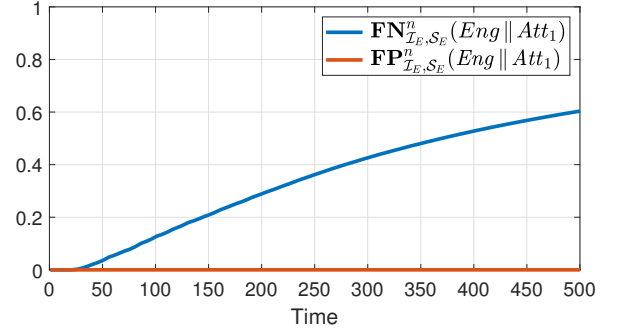


Figure 2. Impact of the first attack on the  $Eng$  system.

Now, we move to our composite system  $Sys$  of Example 3, and suppose that the attack  $Att_1$  tampers with the left engine only. By relying on the compositionality results of our metrics we can infer the following result:

**Proposition 4** (Impact of  $Att_1$  on  $Sys$  when attacking  $Eng^L$ ). *Let  $Sys$  be the system defined in Example 3 with attacker's goal indicators  $\mathcal{I}_{Sys}$  and detection policy  $\mathcal{P}_{Sys}$ , as defined in Example 5. Let  $\mathcal{I} = \mathcal{I}_{Sys}$ ,  $\mathcal{P} = \mathcal{P}_{Sys}$  and  $Att_1^L$  be  $Att_1\{temp^L/temp\}\{cool^L/cool\}$ . Then,*

- $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_1^L) = 0.7 \cdot \mathbf{FN}_{\mathcal{I}_L,\mathcal{P}_L}^n(Eng^L \parallel Att_1^L)$
- $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_1^L) = 0.7 \cdot \mathbf{FP}_{\mathcal{I}_L,\mathcal{P}_L}^n(Eng^L \parallel Att_1^L)$ .

Proposition 4 says that, after 500 time instants, the metric  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$  of the composite system approaches the value  $0.7 \cdot 0.6 = 0.42$ . As expected, the attack does not affect the false positives, i.e.,  $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_1^L) = 0$ .

The *second attack* is an integrity attack to the sensor of the temperature of the  $Eng$  system. The attack adds a negative offset 0.6 to the temperature detected by the controller:

$$Att_2 = \text{rec } X.\text{read } temp(y).\text{write } temp\langle y-0.6 \rangle.\text{tick}.X$$

In this attack, we have a different situation as the IDS raises an alert each time we have  $100 < temp \leq 100.6$  and  $cool = \text{off}$ . As a consequence, this attack is *not stealthy*.

The impact of our second attack in the first 300 time instants is represented in Figure 3. The effects of this attack are recorded in the goal indicator  $stress$ ; the *false negatives* grow with the size of the attack window, reaching the value 0.18 after 300 time instants. As regards the *false positives*, as soon as the attack begins, the metric  $\mathbf{FP}_{\mathcal{I}_E,\mathcal{P}_E}^n$  reaches the value 0.1 and then stabilises. This indicates that, in average, the system under attack experiences 10 false positives every 100 time units, due to unjustified alert signals at channel  $warn$ .

The relatively low number of false negatives is basically due to the mitigation activity of the IDS (consisting in slowing down the engine) that delays the effects of the attack. In fact, without that mitigation the average number of false negatives would have been significantly higher, as depicted in Figure 4. Of course, a smarter attacker that would also drop the controller's command to slow down the  $speed$  (control variable) would cancel the mitigation efforts of the IDS.

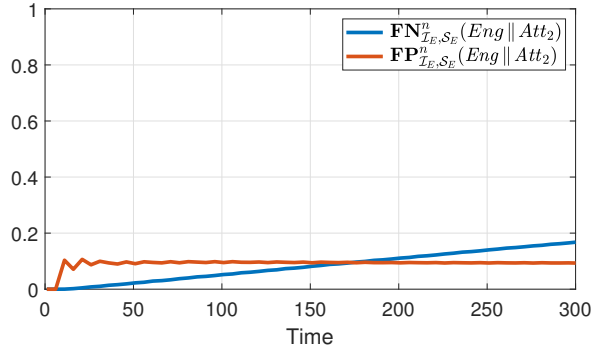


Figure 3. Impact of the second attack on the *Eng* system.

Now, we move to our composite system *Sys* of Example 3, and suppose that the attack *Att*<sub>2</sub> tampers with the right engine only. By relying on the compositionality results of our metrics we can infer the following result:

**Proposition 5** (Impact of *Att*<sub>2</sub> on *Sys* when attacking *Eng*<sup>R</sup>). *Let Sys be the system defined in Example 3 with attacker's goal indicators  $\mathcal{I}_{Sys}$  and detection policy  $\mathcal{P}_{Sys}$ , as defined in Example 5. Let  $\mathcal{I} = \mathcal{I}_{Sys}$  and  $\mathcal{P} = \mathcal{P}_{Sys}$ . Let  $Att_2^R$  be the attack  $Att_2\{^{tempR}/_{temp}\}$ . Then,*

- $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_2^R) = 0.7 \cdot \mathbf{FN}_{\mathcal{I}_R,\mathcal{P}_R}^n(Eng^R \parallel Att_2^R)$
- $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_2^R) = 0.7 \cdot \mathbf{FP}_{\mathcal{I}_R,\mathcal{P}_R}^n(Eng^R \parallel Att_2^R)$ .

The false negatives of the composite system under attack reach the value  $0.7 \cdot 0.18 = 0.126$ . The false positives reach the value  $0.7 \cdot 0.1 = 0.07$ , and then stabilises.

Our *third attack* is a *DoS/integrity* attack on the actuator driving the cooling systems. The attacker tampers with the control variable *cool* to possibly drop two consecutive commands sent by the controller to turn on the cooling system. Once dropped, these commands are replaced with commands to turn off the cooling system. Formally, the attack is the following:

```

Att3 = rec X.tick5.DropForge
DropForge = rec Y.read cool(y).[y = on]
           {write cool(off).OnceMore},
           {write cool(off).tick.Y}
OnceMore = tick.read cool(y).write cool(off).tick.X

```

The attack repeatedly alternates between a *stand-by phase*, lasting 5 time units, and a *drop-and-forge phase*, lasting 2 time units, in which an integrity attack on the variable *cool* may occur once or twice, depending whether the controller transmits one or two (consecutive) commands to turn on the cooling system. Here, the goal indicator *stress* gets increased as the attacker may suddenly turn off the cooling system, reducing the 5-ticks cooling cycle. Thus, either the temperature remains just above the threshold 100 or the temperature drops below 100, but then it rapidly rises above that threshold.

The impact of this attack in the first 300 time instants is represented in Figure 5. The number of undetected and/or underestimated situations of overstress recorded by the goal indicator *stress* rapidly grows with the duration of the attack.

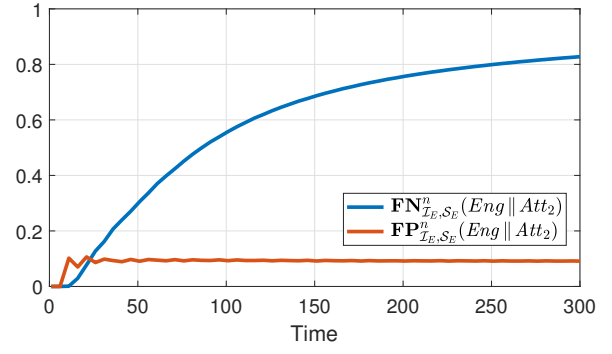


Figure 4. Impact of the second attack with detection but no mitigation.

After only 80 time instants the metric  $\mathbf{FN}_{\mathcal{I}_E,\mathcal{P}_E}^n$  reaches already the value 0.8; after 300 time instants it reaches the value 0.9. As regards the false positives, there is a small peak in the first time instants of malicious activity as the IDS raises alarms which are not related to increased overstress of the system. Thereafter, the average number of false positives quickly drops to 0. From the magnitude of false negatives it follows that the IDS misses most of the attacks targeting the stress of the system; as a consequence, basically, this is a *stealthy attack*.

Now, we move to our composite system *Sys* of Example 3, and suppose that the attack *Att*<sub>3</sub> tampers with the right engine only. By relying on the compositionality results of our metrics we can infer the following result:

**Proposition 6** (Impact of *Att*<sub>3</sub> on *Sys* when attacking *Eng*<sup>R</sup>). *Let Sys be the system defined in Example 3 with attacker's goal indicators  $\mathcal{I}_{Sys}$  and detection policy  $\mathcal{P}_{Sys}$ , as defined in Example 5. Let  $\mathcal{I} = \mathcal{I}_{Sys}$  and  $\mathcal{P} = \mathcal{P}_{Sys}$ . Let  $Att_3^R$  be the attack  $Att_3\{^{coolR}/_{cool}\}$ . Then,*

- $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_3^R) = 0.7 \cdot \mathbf{FN}_{\mathcal{I}_R,\mathcal{P}_R}^n(Eng^R \parallel Att_3^R)$
- $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_3^R) = 0.7 \cdot \mathbf{FP}_{\mathcal{I}_R,\mathcal{P}_R}^n(Eng^R \parallel Att_3^R)$ .

Here, after 300 time instants, the metric  $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n$  of the system *Sys* reaches the value  $0.7 \cdot 0.9 = 0.63$ . As expected, after a short transitory phase, false positives drops to 0.

Finally, we might think of a *fourth attack*  $Att_4 = Att_1^L \parallel Att_2^R$  consisting of two *colluding attacks* focusing on the two engines, respectively: *Att*<sub>1</sub> tampering with the left engine, and *Att*<sub>2</sub> operating on the right engine. In this case, as expected, we would have both false negatives and false positives. However, the attacker's goal indicators of the two engines will let us know where these false alert signals are actually coming from.

**Proposition 7** (Impact of *Att*<sub>4</sub> on *Sys* when attacking *Eng*<sup>L</sup> and *Eng*<sup>R</sup>). *Let Sys be the system defined in Example 3 with attacker's goal indicators  $\mathcal{I}_{Sys}$  and detection policy  $\mathcal{P}_{Sys}$ , as defined in Example 5. Let  $\mathcal{I} = \mathcal{I}_{Sys}$  and  $\mathcal{P} = \mathcal{P}_{Sys}$ . Let  $n = 300$  and  $Att_4$  be the attack  $Att_1^L \parallel Att_2^R$ . Then,*

- $\mathbf{FN}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_4) = \min(1, (0.7 \cdot \mathbf{FN}_{\mathcal{I}_L,\mathcal{P}_L}^n(Eng^L \parallel Att_1^L) + 0.7 \cdot \mathbf{FN}_{\mathcal{I}_R,\mathcal{P}_R}^n(Eng^R \parallel Att_2^R)))$
- $\mathbf{FP}_{\mathcal{I},\mathcal{P}}^n(Sys \parallel Att_4) = \min(1, (0.7 \cdot \mathbf{FP}_{\mathcal{I}_L,\mathcal{P}_L}^n(Eng^L \parallel Att_1^L) + 0.7 \cdot \mathbf{FP}_{\mathcal{I}_R,\mathcal{P}_R}^n(Eng^R \parallel Att_2^R)))$ .

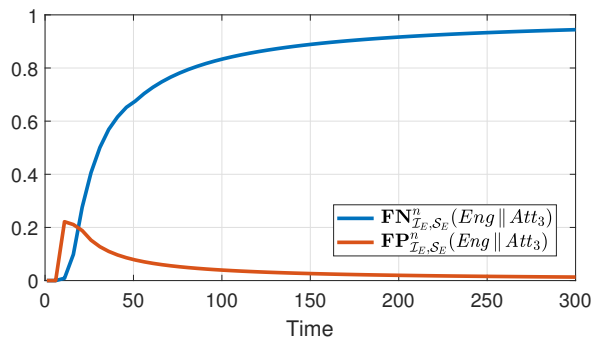


Figure 5. Impact of the third attack on the *Eng* system.

Here, after 300 time instants, the average false negatives approaches the value  $0.7 \cdot 0.5 + 0.7 \cdot 0.18 = 0.476$ , due to the malicious activity on the left engine  $Eng^L$ . Furthermore, the average false positives reaches the value  $0.7 \cdot 0.1 = 0.07$  (and then stabilises), due to the malicious activity on  $Eng^R$ .

## V. CONCLUSIONS AND RELATED WORK

We defined two probabilistic metrics,  $\text{FN}_{\mathcal{I}, \mathcal{P}}^n$  and  $\text{FP}_{\mathcal{I}, \mathcal{P}}^n$ , to estimate the impact of (possibly stealthy) cyber-physical attacks in terms of the average false negatives and the average false positives generated by an arbitrary IDS. Our metrics build on the following concepts: (i) a timed probabilistic labelled transition semantics to formally describe the dynamics of (possibly compromised) CPSs; (ii) a set  $\mathcal{I}$  of weighted attacker’s goal indicators whose runtime values denote how close is the attacker to reach these goals (equipment or production damage, compliance violations); (iii) a detection policy  $\mathcal{P}$  that, given an alert signal (raised by the IDS) and an attacker’s goal indicator in  $\mathcal{I}$ , returns a statically-determined estimate on the progresses of the detected attack in achieving the goal with that indicator.

Our impact metrics can be computed in a compositional manner on those CPSs whose sub-systems may physically interact only if the associated logical components agree on when and how that interaction should happen. Examples of such kind of systems have been discussed in the Introduction.

We showed the computational feasibility of our metrics with the help of a non-trivial use case and four significant cyber-physical attacks, both implemented in MATLAB [15].

As *future work*, we plan to improve the computational efficiency of our metrics by selecting execution traces of the CPS under analysis via Monte Carlo simulation methods [18]. Of course, simulation-based solutions do not ensure results with a 100% confidence. However, the number of executions that our simulator should perform to reach a desired confidence can be determined using well-known statistical techniques [19]. This will allow us to scale our impact metrics on larger CPSs that are not necessarily compositional, according to our definitions.

*Related Work:* Different methodologies for assessing the impact of attacks targeting CPSs have been proposed. Here, we discuss and compare the papers closest to our work.

Urbina et al. [13] addressed the problem of detecting greedy attacks to sensors or actuators, but not both, as we do. Based

on a strong adversary model that is always able to bypass *stateful* IDSs (relying on cumulative sums - CUSUM), they introduced an evaluation metric for attack-detection algorithms that quantifies the negative impact of stealthy attacks and the inherent trade-off with false positives. In practice, in order to define the impact of stealthy attacks the authors select one (or more) variable(s) of interest in the process they want to control. Then, the impact of the undetected attack will measure how much can the attacker drive that value towards its intended goal per unit of time. More precisely, their metric measures the tradeoff between the maximum deviation per time unit imposed by stealthy attacks and the expected time between false positives. Their metric relies exclusively on observable physical variables without explicitly addressing equipment damage, production damage, or compliance violations [6]. We do that by introducing attacker’s goal indicators  $\mathcal{I}$  to capture manipulations on (unobservable) physical variables whose evolution is beyond the control of IDSs.

Umsonst et al. [20] provided the mathematical theory to expand the analysis of [13] based on CUSUM IDSs. In particular, they considered attacks that can tamper with both sensors and actuators at the same time, and can optimise their malicious patterns over a window of time. However, the works of [13] and [20] neglect the influence of the uncertainty of the models and do not propose substitutes for their infinity-norm-based metrics to be used in stochastic systems.

Huang et al. [21] proposed a risk assessment method that uses a Bayesian network to model the attack propagation process and infers the probabilities of sensors and actuators to be compromised. These probabilities are fed into a stochastic hybrid system model to predict the evolution of the physical process being controlled. Then, the security risk is quantified by evaluating the system availability with the model.

More recently, Milošević et al. [22] overcome the limitations of the infinity-norm-based metrics of [13], [20] by proposing two metrics to measure the impact of stealthy cyber-attacks in stochastic linear networked control systems. The first metric returns the probability that some of the critical states leave a safety region; the second metric measures the deviation of such critical states with respect to the expected values. Their metrics relies on two tuning parameters:  $N \in \mathbb{N}$ , denoting the time window over which the impact is estimated (similar to our parameter  $n$ ), and  $\epsilon \geq 0$ , *i.e.*, the stealthiness level based on the Kullback-Leibler divergence. Somehow, these two metrics together provide an information similar to what we obtain with our metric  $\text{FN}_{\mathcal{I}, \mathcal{P}}^n$ , except for the fact that we rely on the detection policy  $\mathcal{P}$  to distill stealthiness instead of using Kullback-Leibler divergence. Unlike us, the authors do not take into consideration false positives introduced by the attacker to mislead system engineers.

Other works on impact evaluation of attacks targeting CPSs are the following. Genge et al. [23] introduced a methodology, inspired by research in system dynamics and sensitivity analysis, to compute the covariances of the observed variables before and after the execution of specific interventions involving the control variables. Bilis et al. [24] proposed five metrics derived

from network theory to assess the impacts of cyber attacks on power systems. Sgouras et al. [25] evaluated the impact of attacks on a simulated smart metering infrastructure. Sridhar and Govindarasu [26] evaluated the impact of attacks on wide-area frequency control applications in power systems.

As regards *formal methodologies*, Lanotte et al. [11], [10] defined hybrid process calculi to model both CPSs and cyber-physical attacks; our calculus is a probabilistic simplification of theirs. The paper [10] proposes a formalisation of physical impact of an attack in terms of the *potential perturbation* that might be introduced by the attacker during a compromised execution of the system. However, the model in that paper does not support probabilities to represent system uncertainty. Thus, the notion of impact in [10] does not take into account the attacker's chances of achieving an attack with that physical impact/damage: it just looks at the existence of a malicious execution trace reaching that impact/damage; it says nothing about the likelihood of such an execution actually taking place. More recently, Lanotte et al. [27] have used a discrete-time generalisation of Desharnais et al.'s *weak bisimulation metric* [28] to estimate the impact of attacks targeting sensor devices of IoT systems; however, their impact metric takes into consideration only the logical effects on (communication channels of) the IoT system under attack, but not the physical ones. The objective of this paper is to bring the ideas from these two papers together by defining impact metrics that take into account both the physical impact and the likelihood to successfully completing an attack with that physical impact.

Finally, we wish to underline that *formal methodologies* are increasingly used to lay theoretical foundations to reason about cyber-physical systems and their vulnerabilities (see, e.g., [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39]).

#### ACKNOWLEDGMENT

We thank the anonymous reviewers for their comments. Massimo Merro has been partially supported by the project "Dipartimenti di Eccellenza 2018–2022", funded by the Italian Ministry of Universities and Research (MUR).

#### REFERENCES

- [1] R. Rajkumar, I. Lee, L. Sha, and J. A. Stankovic, "Cyber-physical systems: the next computing revolution," in *DAC*. ACM, 2010, pp. 731–736.
- [2] Y. Huang, A. A. Cárdenas, S. Amin, Z. Lin, H. Tsai, and S. Sastry, "Understanding the physical and economic consequences of attacks on control systems," *IJCIP*, vol. 2, no. 3, pp. 73–83, 2009.
- [3] N. Falliere, L. Murchu, and E. Chien, "W32.Stuxnet Dossier," 2011.
- [4] J. Slowik, "Anatomy of an attack: Detecting and defeating CRASHOVER-RIDE," *VB2018, October*, 2018.
- [5] A. Di Pinto, Y. Dragoni, and A. Carcano, "TRITON: The First ICS Cyber Attack on Safety Instrument Systems," in *Black Hat USA*, 2018.
- [6] D. Gollmann, P. Gurikov, A. Isakov, M. Krotofil, J. Larsen, and A. Winnicki, "Cyber-Physical Systems Security: Experimental Analysis of a Vinyl Acetate Monomer Plant," in *CPSS*. ACM, 2015, pp. 1–12.
- [7] R. Milner, J. Parrow, and D. Walker, "A Calculus of Mobile Processes," *Inf. Comput.*, vol. 100, no. 1, pp. 1–40, 1992.
- [8] L. Cardelli and A. Gordon, "Mobile Ambients," *Theoretical Computer Science*, vol. 240, no. 1, pp. 177–213, 2000.
- [9] M. Abadi, B. Blanchet, and C. Fournet, "The applied pi calculus: mobile values, new names and secure communication," *JACM*, vol. 65:1-41, 2018.
- [10] R. Lanotte, M. Merro, A. Munteanu, and L. Viganò, "A formal approach to physics-based attacks in cyber-physical systems," *ACM Trans. Priv. Secur.*, vol. 23, no. 1, pp. 3:1–3:41, 2020.
- [11] R. Lanotte, M. Merro, and S. Tini, "A probabilistic calculus of cyber-physical systems," *Information and Computation*, no. 104618, 2020.
- [12] R. Segala, "Modeling and Verification of Randomized Distributed Real-Time Systems," Ph.D. dissertation, MIT, 1995.
- [13] D. I. Urbina, J. A. Giraldo, A. A. Cárdenas, N. O. Tippenhauer, J. Valente, M. A. Faisal, J. Ruths, R. Candell, and H. Sandberg, "Limiting the Impact of Stealthy Attacks on Industrial Control Systems," in *CCS*. ACM, 2016, pp. 1092–1105.
- [14] A. A. Cárdenas, J. S. Baras, and K. Seamon, "A Framework for the Evaluation of Intrusion Detection Systems," in *IEEE Symposium on Security and Privacy*. IEEE Computer Society, 2006, pp. 63–77.
- [15] MATLAB, 9.7.0.1190202 (R2019b). US: The MathWorks Inc., 2018.
- [16] M. Hennessy and T. Regan, "A Process Algebra for Timed Systems," *Information and Computation*, vol. 117, no. 2, pp. 221–239, 1995.
- [17] R. M. Keller, "Formal verification of parallel programs," *Communications of the ACM*, vol. 19, pp. 371–384, 1976.
- [18] W. K. Hastings, "Monte carlo sampling methods using markov chains and their applications," *Biometrika*, vol. 57, no. 1, pp. 97–109, 1970.
- [19] H. Chernoff, "A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations," *The Annals of Mathematical Statistics*, vol. 23, no. 4, pp. 493–507, 12 1952.
- [20] D. Umsonst, H. Sandberg, and A. A. Cárdenas, "Security analysis of control system anomaly detectors," in *ACC*. IEEE, 2017, pp. 5500–5506.
- [21] K. Huang, C. Zhou, Y. Tian, S. Yang, and Y. Qin, "Assessing the Physical Impact of Cyberattacks on Industrial Cyber-Physical Systems," *IEEE Trans. Industrial Electronics*, vol. 65, no. 10, pp. 8153–8162, 2018.
- [22] J. Milošević, H. Sandberg, and K. H. Johansson, "Estimating the Impact of Cyber-Attack Strategies for Stochastic Networked Control Systems," *IEEE Trans. Control. Netw. Syst.*, vol. 7, no. 2, pp. 747–757, 2020.
- [23] B. Genge, I. Kiss, and P. Haller, "A system dynamics approach for assessing the impact of cyber attacks on critical infrastructures," *Int. J. Critical Infrastructure Protection*, vol. 10, pp. 3–17, 2015.
- [24] E. I. Bilis, W. Kröger, and N. Cen, "Performance of Electric Power Systems Under Physical Malicious Attacks," *IEEE Systems Journal*, vol. 7, no. 4, pp. 854–865, 2013.
- [25] K. I. Sgouras, A. I. Birda, and D. L. Labridis, "Cyber attack impact on critical Smart Grid infrastructures," in *PES ISGT*. IEEE, 2014, pp. 1–5.
- [26] S. Sridhar and M. Govindarasu, "Model-Based Attack Detection and Mitigation for Automatic Generation Control," *IEEE Trans. Smart Grid*, vol. 5, no. 2, pp. 580–591, 2014.
- [27] R. Lanotte, M. Merro, and S. Tini, "Towards a formal notion of impact metric for cyber-physical attacks," in *iFM*, ser. LNCS, vol. 11023. Springer, 2018, pp. 296–315.
- [28] J. Desharnais, R. Jagadeesan, V. Gupta, and P. Panangaden, "The Metric Analogue of Weak Bisimulation for Probabilistic Processes," in *LICS*. IEEE Computer Society, 2002, pp. 413–422.
- [29] V. Nigam, C. Talcott, and A. Urquiza, "Towards the Automated Verification of Cyber-Physical Security Protocols: Bounding the Number of Timed Intruders," in *ESORICS*, ser. LNCS, vol. 9879:450-470, 2016.
- [30] M. Rocchetto and N. Tippenhauer, "On attacker models and profiles for cyber-physical systems" in *ESORICS*, ser. LNCS, vol. 9879:427-449, 2016.
- [31] B. Bohrer and A. Platzer, "A Hybrid, Dynamic Logic for Hybrid-Dynamic Information Flow," in *LICS*. ACM/IEEE, 2018, pp. 115–124.
- [32] V. Nigam and C. L. Talcott, "Formal Security Verification of Industry 4.0 Applications," in *ETFA*. IEEE, 2019, pp. 1043–1050.
- [33] C. Bodei, S. Chessa, and L. Galletta, "Measuring security in IoT communications," *Theor. Comput. Sci.*, vol. 764, pp. 100–124, 2019.
- [34] M. Barrère, C. Hankin, N. Nicolaou, D. G. Eliades, and T. Parisini, "Measuring cyber-physical security in industrial control systems via minimum-effort attack strategies," *J. Inf. Secur. Appl.*, vol. 52, 2020.
- [35] L. Lin, Y. Zhu, and R. Su, "Synthesis of covert actuator attackers for free," *Discret. Event Dyn. Syst.*, vol. 30, no. 4, pp. 561–577, 2020.
- [36] R. Lanotte, M. Merro, and A. Munteanu, "Runtime Enforcement for Control System Security," in *CSF*. IEEE, 2020, pp. 246–261.
- [37] C. Bernardeschi, A. Domenici, and M. Palmieri, "Formalization and co-simulation of attacks on cyber-physical systems," *J. Comput. Virol. Hacking Tech.*, vol. 16, no. 1, pp. 63–77, 2020.
- [38] D. Bresolin, P. Collins, L. Geretti, R. Segala, T. Villa, and S. Z. Gonzalez, "A computable and compositional semantics for hybrid automata," in *HSCC*. ACM, 2020, pp. 18:1–18:11.
- [39] I. Jahandideh, F. Ghassemi, and M. Sirjani, "An actor-based framework for asynchronous event-based cyber-physical systems," *Software and Systems Modeling*, pp. 1–25, in press.

## A. Proofs of Section III

*Proof of Proposition 1.* Given a trace  $M_1 \xrightarrow{\alpha_1 \dots \alpha_n} M_{n+1}$ , with  $M_1 = Eng$ , for  $j \in 1..n+1$ , we let  $t_j$ ,  $c_j$ ,  $sp_j$  and  $s_j$  denote the value of the variables *temp*, *cool*, *speed* and *stress*, respectively, in the physical state of  $M_j$ . We recall that  $\mathcal{I}_E^- = \{stress\}$ ,  $\mathcal{P}_E(\overline{warn}\langle ID, hot \rangle, stress) = 1$  and  $\mathcal{P}_E(\alpha, stress) = 0$ , for all  $\alpha \neq \overline{warn}\langle ID, hot \rangle$ . Thus, in order to derive the result it is sufficient to show that any trace  $M_1 \xrightarrow{\alpha_1 \dots \alpha_{n-1}} M_n$  satisfies: 1)  $\alpha_j \neq \overline{warn}\langle ID, hot \rangle$ , for  $j \in 1..n-1$ ; and 2)  $s_j = 0$ , for  $j \in 1..n$ . Indeed,  $\mathbf{FP}_{\mathcal{I}_E^-, \mathcal{P}_E}^n(Eng) = 0$  follows by 1), and  $\mathbf{FN}_{\mathcal{I}_E^-, \mathcal{P}_E}^n(Eng) = 0$  follows by 2). We prove these two properties separately.

*Case 1).* Since  $\overline{warn}\langle ID, hot \rangle$  can be emitted uniquely by the process *IDS* when  $(z_{temp} > 100 \wedge z_{cool} = \text{off})$  is true, it is sufficient to prove that the condition  $(t_j > 100 \wedge c_j = \text{off})$  is never satisfied, for any  $j \in 1..n$ . We know that  $t_1 = 95$  and  $c_1 = \text{off}$ . Moreover, the temperature rises by at most 1.2 degrees per time unit. Therefore, at some instant  $k$  the temperature exceeds the threshold, i.e.,  $100 < t_k < 101.2$ . The controller *Ctrl* immediately turns on the cooling system, which remains active for 5 time instants, namely  $c_j = \text{on}$  for  $j \in k..k+4$ . This ensures that  $(t_j > 100 \wedge c_j = \text{off})$  remains false for  $j \in k..k+4$ . In this time interval, the cooling system lowers the temperature by at least  $0.8 \cdot 5 = 4$  degrees. As a consequence, we have  $t_{k+5} < 97.2$ . This implies that  $(t_j > 100 \wedge c_j = \text{off})$  remains false at instant  $k+5$ . The reasoning can now be reiterated to prove the required result.

*Case 2).* We will show that for any  $j \geq 0$  we have: (a) if  $t_j > 100$  then  $t_{j+1} \leq 100$  or  $t_{j+2} \leq 100$ ; (b) if  $t_j > 100$  and  $t_{j+1} \leq 100$ , then  $t_{j+1+k} \leq 100$  for  $k \in 1 \dots 4$ . From (a) and (b) it follows that for any arbitrary time interval  $i..i+5$  at most two values in  $\{t_i, \dots, t_{i+5}\}$  may be above the threshold 100. This ensures that *stress* is never incremented and maintains the initial value 0. In order to show (a), we note that the temperature cannot rise above  $100 + 1.2 = 101.2$  degrees. Hence  $t_j \in (100, 101.2]$ . The cooling system is turned on implying  $t_{j+1} \in (100 - 1.2, 101.2 - 0.8] = (98.8, 100.4]$  and  $t_{j+2} \in (98.8 - 1.2, 100.4 - 0.8] = (97.6, 99.6]$ . This gives (a). In order to show (b), from  $t_j > 100$  and  $t_{j+1} \leq 100$  we infer  $t_{j+1} \in (100 - 1.2, 100] = (98.8, 100]$ . The cooling system operates for further 4 time units and  $t_{j+1+k} \in (98.8 - k \cdot 1.2, 100 - k \cdot 0.8]$ , namely  $t_{j+1+k} \leq 100$ , for  $k \in 1..4$ .  $\square$

In order to prove Theorem 1, we need a preliminary result to formalise the relationship between the attackers' progresses associated with the traces of a CPS  $M_1 \uplus M_2$ , and those associated with the traces of the two sub-systems  $M_1$  and  $M_2$ .

**Lemma 1.** *Let  $M_1$  and  $M_2$  be two CPSs with disjoint sets of weighted attacker's goal indicators  $\mathcal{I}_1 = \{(i_1^1, r_1^1), \dots, (i_{k_1}^1, r_{k_1}^1)\}$  and  $\mathcal{I}_2 = \{(i_1^2, r_1^2), \dots, (i_{k_2}^2, r_{k_2}^2)\}$ , and detection policies  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , respectively. Assume that for any  $i, j \in \{1, 2\}$ ,  $i \neq j$ , whenever  $\bar{c}m \in \text{out}(M_i)$  and  $cm \in \text{inp}(M_j)$  then  $\mathcal{P}_i(\bar{c}m, h) = 0$ , for any  $h \in \mathcal{I}_i^-$ . Then,*

for two tuples of reals  $\bar{q}^1 = q_1^1, \dots, q_{k_1}^1$  and  $\bar{q}^2 = q_1^2, \dots, q_{k_2}^2$ , the following two properties are equivalent:

- 1)  $M_1 \uplus M_2 \xrightarrow{\bar{q}^1, \bar{q}^2} \mathcal{P} \gamma$ , with  $\mathcal{P} = p_1 \cdot \mathcal{P}_1 \uplus p_2 \cdot \mathcal{P}_2$ ;
- 2)  $M_1 \xrightarrow{\bar{q}^1} \mathcal{P}_1 \gamma_1$ ,  $M_2 \xrightarrow{\bar{q}^2} \mathcal{P}_2 \gamma_2$  and  $\gamma = \gamma_1 \uplus \gamma_2$ ,

where  $\gamma_1 \uplus \gamma_2$  is the distribution over CPSs such that  $(\gamma_1 \uplus \gamma_2)(N_1 \uplus N_2) = \gamma_1(N_1) \cdot \gamma_2(N_2)$ , for all CPSs  $N_1$  and  $N_2$ .

*Proof:* We show that property 1) implies 2), the other implication is similar. Let  $M_1 = \langle \varepsilon_1; \xi_{V_1} \rangle \bowtie P_1$  and  $M_2 = \langle \varepsilon_2; \xi_{V_2} \rangle \bowtie P_2$ , for evolution laws  $\varepsilon_1$  and  $\varepsilon_2$ , state functions  $\xi_{V_1}$  and  $\xi_{V_2}$  and processes  $P_1$  and  $P_2$ . Let  $N^0 = M_1 \uplus M_2$ . By Definition 9, the transition in 1) entails:

- $N^0 \xrightarrow{\alpha_1} \dots \xrightarrow{\alpha_n} N^n$ , where  $\alpha_j \neq \text{tick}$ , for  $j \in 1..n$
- $N^n \xrightarrow{\text{tick}} \gamma$ .

All systems  $N^j$ , with  $j \in 0..n$ , have the form  $N^j = \langle \varepsilon_1 \uplus \varepsilon_2; \xi_{V_1} \uplus \xi_{V_2} \rangle \bowtie P_1^j \parallel P_2^j$ , for some processes  $P_1^j$  and  $P_2^j$ . The proof of this statement follows by induction on  $j$  and by inspection on the transition rules of Tables I and II. Then, the transitions  $N^j \xrightarrow{\alpha_{j+1}} N^{j+1}$ , for  $j \in 0..n-1$ , can be derived only because  $P_1^j \parallel P_2^j \xrightarrow{\beta_{j+1}} P_1^{j+1} \parallel P_2^{j+1}$ , for some  $\beta_{j+1}$ .

The transitions  $N^n \xrightarrow{\text{tick}} \gamma$  can be derived only because  $P_1^n \parallel P_2^n \xrightarrow{\text{tick}} P_1^{n+1} \parallel P_2^{n+1}$  by the application of either rule (ReadIDS) or rule (Time). Thus, the distribution  $\gamma$  of the transition  $N^n \xrightarrow{\text{tick}} \gamma$  has the form  $\gamma = \gamma_1 \uplus \gamma_2$ , where  $\gamma_1 = \varepsilon_1(\xi_{V_1}) \bowtie P_1^{n+1}$  and  $\gamma_2 = \varepsilon_2(\xi_{V_2}) \bowtie P_2^{n+1}$ , for some processes  $P_1^{n+1}$  and  $P_2^{n+1}$ .

Thus, let  $N_1^j = \langle \varepsilon_1; \xi_{V_1} \rangle \bowtie P_1^j$  and  $N_2^j = \langle \varepsilon_2; \xi_{V_2} \rangle \bowtie P_2^j$ , for  $j \in 1..n$ , the transitions of CPSs  $N^j$ , for  $j \in 0..n$ , follow from transitions of the associated processes. There are 4 cases.

*Case (i).*  $P_1^j \parallel P_2^j \xrightarrow{\beta_{j+1}} P_1^{j+1} \parallel P_2^{j+1}$  is derived by an application of rule (Par) in Table I because  $P_1^j \xrightarrow{\beta_{j+1}} P_1^{j+1}$  and  $P_2^{j+1} = P_2^j$ . Now, whatever is the rule in Table II used to derive  $N^j \xrightarrow{\alpha_{j+1}} N^{j+1}$  from  $P_1^j \parallel P_2^j \xrightarrow{\beta_{j+1}} P_1^{j+1} \parallel P_2^{j+1}$ , this rule allows us to derive  $N_1^j \xrightarrow{\alpha_{j+1}} N_1^{j+1}$  as well, from  $P_1^j \xrightarrow{\beta_{j+1}} P_1^{j+1}$ . As  $\mathcal{P} = p_1 \cdot \mathcal{P}_1 \uplus p_2 \cdot \mathcal{P}_2$ , we have  $\mathcal{P}_1(\alpha_{j+1}, i_l^1) = \mathcal{P}(\alpha_{j+1}, i_l^1)$ , for  $\alpha_{j+1} \in \text{out}(M_1)$  and  $l = 1..k_1$ .

*Case (ii).*  $P_1^j \parallel P_2^j \xrightarrow{\beta_{j+1}} P_1^{j+1} \parallel P_2^{j+1}$  is derived by applying rule (Par) because  $P_2^j \xrightarrow{\beta_{j+1}} P_2^{j+1}$ . Similarly to case (i) we have  $\mathcal{P}_2(\alpha_{j+1}, i_l^2) = \mathcal{P}(\alpha_{j+1}, i_l^2)$ , for  $\alpha_{j+1} \in \text{out}(M_2)$  and  $l = 1..k_2$ .

*Case (iii).*  $P_1^j \parallel P_2^j \xrightarrow{\beta_{j+1}} P_1^{j+1} \parallel P_2^{j+1}$  is derived by rule (Comm) in Table I because  $P_1^j \xrightarrow{\beta_{j+1}^1} P_1^{j+1}$  and  $P_2^j \xrightarrow{\beta_{j+1}^2} P_2^{j+1}$ . Without loss of generality, we assume  $\beta_{j+1}^1 = \bar{c}m$  and  $\beta_{j+1}^2 = cm$ . By rules (Receive) and (Send), we infer  $N_1^j \xrightarrow{\alpha_{j+1}^1} N_1^{j+1}$  and  $N_2^j \xrightarrow{\alpha_{j+1}^2} N_2^{j+1}$ , with  $\alpha_{j+1}^1 = \beta_{j+1}^1$  and  $\alpha_{j+1}^2 = \beta_{j+1}^2$ . Since  $\alpha_{j+1}^1 = \bar{c}m \in \text{out}(M_1)$  and  $\alpha_{j+1}^2 = cm \in \text{inp}(M_2)$ , by applying the hypothesis of the lemma we have  $\mathcal{P}_1(\alpha_{j+1}^1, i_l^1) = 0$ . Thus,  $\alpha_{j+1}^2 \notin \text{out}(M_2)$  and  $\alpha_{j+1} = \tau \notin \text{out}(M)$ .

*Case (iv).*  $P_1^n \parallel P_2^n \xrightarrow{\text{tick}} P_1^{n+1} \parallel P_2^{n+1}$  is derived by an application of rule (Timepar) in Table I because  $P_1^n \xrightarrow{\text{tick}} P_1^{n+1}$  and  $P_2^n \xrightarrow{\text{tick}} P_2^{n+1}$ . In this case, we have:  $N_1^n \xrightarrow{\text{tick}} N_1^{n+1}$

$\gamma_1$  with  $\gamma_1 = \varepsilon_1(\xi_{V_1}) \bowtie P_1^{n+1}$ , and  $N_2^n \xrightarrow{\text{tick}} \gamma_2$  with  $\gamma_2 = \varepsilon_2(\xi_{V_2}) \bowtie P_2^{n+1}$ .

Now, from the form of the transitions originating from CPSs  $N_i^j$  and  $N_i^n$ , for  $j \in 0..n-1$  and  $i \in 1..2$ , we derive the following five facts:

- (a)  $M_1 = N_1^0 \xrightarrow{\delta_{1,1}} \dots \xrightarrow{\delta_{1,n_1}} N_1^{n_1} = N_1^n \xrightarrow{\text{tick}} \gamma_1$ ;
- (b) for all  $j \in 1..n_1$  and  $l \in 1..r_1$ , either  $\mathcal{P}_1(\delta_{1,j}, i_l^1) = \mathcal{P}(\alpha_j, i_l^1)$  or both  $\delta_{1,j} \notin \text{out}(M_1)$  and  $\alpha_j \notin \text{out}(M)$ ;
- (c)  $M_2 = N_2^0 \xrightarrow{\delta_{2,1}} \dots \xrightarrow{\delta_{2,n_2}} N_2^{n_2} = N_2^n \xrightarrow{\text{tick}} \gamma_2$ ;
- (d) for all  $j \in 1..n_2$  and  $l \in 1..r_2$ , either  $\mathcal{P}_2(\delta_{2,j}, i_l^2) = \mathcal{P}(\alpha_j, i_l^2)$  or both  $\delta_{2,j} \notin \text{out}(M_2)$  and  $\alpha_j \notin \text{out}(M)$ ;
- (e)  $\delta_{1,j_1} \neq \text{tick}$  and  $\delta_{2,j_2} \neq \text{tick}$ , for  $j_1 \in 1..n_1$ ,  $j_2 \in 1..n_2$ .

From (a) and (e), by applying Definition 9 we derive the transition  $M_1 \xrightarrow{r_1^1, \dots, r_{k_1}^1} \mathcal{P}_1 \gamma_1$ , with  $r_l^1 = \min(1, \sum_{j=1}^{n_1} r_l^{1,j})$ ; here,  $r_l^{1,j} = \mathcal{P}_1(\delta_{1,j}, i_l^1)$  if  $\delta_{1,j} \in \text{out}(M_1)$ , and  $r_l^{1,j} = 0$  otherwise. We recall that by Definition 9, the transition in item 1) of the hypothesis of the lemma entails  $q_l^1 = \min(1, \sum_{j=1}^n q_l^{1,j})$ , where  $q_l^{1,j} = \mathcal{P}(\alpha_j, i_l^1)$  if  $\alpha_j \in \text{out}(M_1)$ , and  $q_l^{1,j} = 0$  otherwise. Hence, by fact (b) it follows that  $r_l^1 = q_l^1$ . This implies  $M_1 \xrightarrow{\bar{q}_1} \mathcal{P}_1 \gamma_1$ , for  $\bar{q}_1 = q_1^1, \dots, q_{k_1}^1$ . By applying a similar reasoning, from facts (c), (d) and (e) we derive  $M_2 \xrightarrow{\bar{q}_2} \mathcal{P}_2 \gamma_2$ . The result follows because  $\gamma = \gamma_1 \uplus \gamma_2$ . ■

Now, everything is in place to prove Theorem 1.

*Proof of Theorem 1.* We give the proof for the metric **FP**. The proof for the metric **FN** is similar. Moreover, we consider the case  $k = 2$ , the extension to  $k > 2$  is straightforward. We reason by induction on  $n$ . The base case,  $n = 0$ , is immediate. We move to the inductive case,  $n > 0$ . Let  $M = M_1 \uplus M_2$ . Assume  $\mathcal{I}_1 = \{(i_1^1, r_1^1), \dots, (i_{k_1}^1, r_{k_1}^1)\}$  and  $\mathcal{I}_2 = \{(i_1^2, r_1^2), \dots, (i_{k_2}^2, r_{k_2}^2)\}$ . Let  $\bar{q}_1$  denote a tuple of reals  $q_1^1, \dots, q_{k_1}^1$  and  $\bar{q}_2$  a tuple of reals  $q_1^2, \dots, q_{k_2}^2$ . Then, we have:

$$\begin{aligned} & \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M) \\ = & \min \left( 1, \max_{M \xrightarrow{\bar{q}_1, \bar{q}_2} \mathcal{P} \gamma} \frac{\sum_{j=1}^{k_1} p_1 \cdot r_j^1 \cdot \max(0, q_j^1 - i_j^1(M)) + \sum_{j=1}^{k_2} p_2 \cdot r_j^2 \cdot \max(0, q_j^2 - i_j^2(M)) + (n-1) \sum_{M' \in \text{supp}(\gamma)} \gamma(M') \cdot \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^{n-1}(M')}{n} \right) \\ = & \min \left( 1, \max_{\substack{M_1 \xrightarrow{\bar{q}_1} \mathcal{P}_1 \gamma_1 \\ M_2 \xrightarrow{\bar{q}_2} \mathcal{P}_2 \gamma_2}} \frac{\sum_{j=1}^{k_1} p_1 \cdot r_j^1 \cdot \max(0, q_j^1 - i_j^1(M)) + \sum_{j=1}^{k_2} p_2 \cdot r_j^2 \cdot \max(0, q_j^2 - i_j^2(M)) + (n-1) \sum_{M'_1 \in \text{supp}(\gamma_1)} \gamma_1(M'_1) \cdot \gamma_2(M'_2) \cdot \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^{n-1}(M'_1 \uplus M'_2)}{n} \right) \\ = & \min \left( 1, \max_{\substack{M_1 \xrightarrow{\bar{q}_1} \mathcal{P}_1 \gamma_1 \\ M_2 \xrightarrow{\bar{q}_2} \mathcal{P}_2 \gamma_2}} \frac{\sum_{j=1}^{k_1} p_1 \cdot r_j^1 \cdot \max(0, q_j^1 - i_j^1(M)) + \sum_{j=1}^{k_2} p_2 \cdot r_j^2 \cdot \max(0, q_j^2 - i_j^2(M)) + (n-1) \sum_{M'_1 \in \text{supp}(\gamma_1)} \gamma_1(M'_1) \cdot \gamma_2(M'_2) \cdot \min(1, p_1 \cdot \mathbf{FP}_{\mathcal{I}_1, \mathcal{P}_1}^{n-1}(M'_1) + p_2 \cdot \mathbf{FP}_{\mathcal{I}_2, \mathcal{P}_2}^{n-1}(M'_2))}{n} \right) \end{aligned}$$

$$\begin{aligned} & p_1 \cdot \sum_{j=1}^{k_1} r_j^1 \cdot \max(0, q_j^1 - i_j^1(M_1)) + \\ & (n-1) \sum_{M'_1 \in \text{supp}(\gamma_1)} \gamma_1(M'_1) p_1 \mathbf{FP}_{\mathcal{I}_1, \mathcal{P}_1}^{n-1}(M'_1) + \\ & p_2 \cdot \sum_{j=1}^{k_2} r_j^2 \cdot \max(0, q_j^2 - i_j^2(M_2)) + \\ & (n-1) \sum_{M'_2 \in \text{supp}(\gamma_2)} \gamma_2(M'_2) p_2 \mathbf{FP}_{\mathcal{I}_2, \mathcal{P}_2}^{n-1}(M'_2) \\ = & \min \left( 1, \max_{\substack{M_1 \xrightarrow{\bar{q}_1} \mathcal{P}_1 \gamma_1 \\ M_2 \xrightarrow{\bar{q}_2} \mathcal{P}_2 \gamma_2}} \frac{\sum_{j=1}^{k_1} p_1 \cdot r_j^1 \cdot \max(0, q_j^1 - i_j^1(M_1)) + (n-1) \sum_{M'_1 \in \text{supp}(\gamma_1)} \gamma_1(M'_1) p_1 \mathbf{FP}_{\mathcal{I}_1, \mathcal{P}_1}^{n-1}(M'_1) + p_2 \cdot \sum_{j=1}^{k_2} r_j^2 \cdot \max(0, q_j^2 - i_j^2(M_2)) + (n-1) \sum_{M'_2 \in \text{supp}(\gamma_2)} \gamma_2(M'_2) p_2 \mathbf{FP}_{\mathcal{I}_2, \mathcal{P}_2}^{n-1}(M'_2)}{n} \right) \\ = & \min(1, p_1 \cdot \mathbf{FP}_{\mathcal{I}_1, \mathcal{P}_1}^n(M_1) + p_2 \cdot \mathbf{FP}_{\mathcal{I}_2, \mathcal{P}_2}^n(M_2)). \end{aligned}$$

Here, step 1 follows by the definition of  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M)$ ; step 2 by an application of Lemma 1; step 3 by inductive hypothesis; step 4 by immediate calculation; last step follows by the definitions of  $\mathbf{FP}_{\mathcal{I}_1, \mathcal{P}_1}^n(M_1)$  and  $\mathbf{FP}_{\mathcal{I}_2, \mathcal{P}_2}^n(M_2)$ . □

It remains to prove both Proposition 2 and Theorem 2. We focus on the proof of Theorem 2; the proof of Proposition 2 is simpler. First, let us introduce some technical notations.

**Notation 4.** Given a set of actions  $\mathcal{A} \subseteq \mathbf{A}$ , we write  $\text{chn}(\mathcal{A})$  to denote the set  $\{c: cm \in \mathcal{A} \text{ or } \bar{c}m \in \mathcal{A}\}$  of the channels used in  $\mathcal{A}$ . Given an execution trace  $t$ , we write  $\text{restr}_{\mathcal{A}}(t)$  to denote the trace resulting by removing in  $t$  all actions which are not in  $\mathcal{A}$ . A set of actions  $\mathcal{A} \subseteq \text{out}(M)$  will be called a set of alert signals for the CPS  $M$  if it contains all and only output actions of  $M$  (and its derivatives) on channels which are never used in input; formally,  $\alpha \in \mathcal{A}$  iff  $\alpha \in \text{out}(M)$  and  $\text{chn}(\alpha) \not\subseteq \text{chn}(\text{inp}(M))$ . We say that a channel  $c$  is fresh with respect to a CPS  $M$  if  $c$  is never used in  $M$ . Finally, given a CPS  $M$  we define the set of tick-derivatives of  $M$  as  $\text{der}_{\text{tick}}(M) = \{M' : M \xrightarrow{t} \text{tick} M'\} \cup \{M\}$ .

Now, we need a definition to formalise how to compare two detection policies in term of their aggressiveness.

Let us consider a CPS  $M$  and a supervisor process  $Sup$ . Assume two sets of alert signals  $\mathcal{A} \subseteq \text{out}(M)$  and  $\mathcal{B} \subseteq \text{out}(Sup)$  such that the signals in  $\mathcal{A}$  are used for communications from  $M$  to the supervisor, while those in  $\mathcal{B}$  are used for communications from the supervisor to the external observer. We also assume that the supervisor  $Sup$  collects all alert signals of  $\mathcal{A}$  and converts them in alert signals of  $\mathcal{B}$ , on communication channels fresh with respect to  $M$ . More formally,

- for all  $\widehat{Sup} \in \text{der}_{\text{tick}}(Sup)$ , whenever  $\widehat{Sup} \xrightarrow{t} \text{tick} \widehat{Sup}$ , with  $\text{tick} \notin t$ , then  $t = c_1 m_1 \dots c_h m_h \bar{d}_1 m'_1 \dots \bar{d}_l m'_l$ , for  $\bar{c}_j m_j \in \mathcal{A}$  and  $\bar{d}_k m'_k \in \mathcal{B}$ ; we call this property *signal conversion*;
- all channels in  $\text{chn}(\mathcal{B})$  are fresh with respect to  $M$ .

Now, if  $Sup$  is able convert all alert signals in  $\mathcal{A} \subseteq \text{out}(M)$  to fresh alert signals in  $\mathcal{B} \subseteq \text{out}(Sup)$  then we can compare the detection policies of the CPS  $M$  and the supervised CPS  $(M \parallel Sup) \setminus \text{chn}(\mathcal{A})$ , in terms of their aggressiveness. Formally,

**Definition 12** (Detection policy comparison). *Let  $M \in \mathcal{M}$  be a CPS with attacker's goal indicators  $\mathcal{I}$ , alert signals  $\mathcal{A} \subseteq \text{out}(M)$  and detection policy  $\mathcal{P}$ . Let  $Sup$  be a pure-logical process, with detection policy  $\mathcal{P}_{Sup}$ , and collecting all alert*

signals  $\mathcal{A}$  converting them in alert signals  $\mathcal{B} \subseteq \text{out}(Sup)$ , on fresh communication channels. Let  $\sim$  be a binary operator to compare real numbers:  $\sim \in \{\leq, =, \geq\}$ .

We write  $\mathcal{P}_{Sup} \sim \mathcal{P}$  when for all  $(M' \parallel Sup') \setminus \text{chn}(\mathcal{A}) \in \text{der}_{\text{tick}}((M \parallel Sup) \setminus \text{chn}(\mathcal{A}))$  and tick-derivatives  $M''$  of  $M'$ , for which  $M' \xrightarrow{t} \xrightarrow{\text{tick}} M''$ , with  $\text{tick} \notin t$ , there is a trace  $Sup' \xrightarrow{t'} \xrightarrow{\text{tick}} Sup''$ ,  $\text{tick} \notin t'$ , such that:

- $\text{restr}_{\mathcal{A}}(t) = \overline{c_1}m_1 \cdots \overline{c_h}m_h$ , for some  $h \in \mathbb{N}$
- $t' = c_1m_1 \cdots c_hm_h \cdot \overline{d_1}m'_1 \cdots \overline{d_l}m'_l$ , for some  $l \in \mathbb{N}$
- $\min(1, \sum_{j=1}^l \mathcal{P}_{Sup}(\overline{d_j}m'_j, i)) \sim \min(1, \sum_{j=1}^h \mathcal{P}(\overline{c_j}m_j, i))$  for all attacker's goal indicators  $i \in \mathcal{I}$ .

In order to prove Theorem 2, we need a preliminary result, stating the relationships between the attackers' progresses associated with the traces of a CPS  $M$  and those associated with the traces of  $(M \parallel Sup) \setminus \text{chn}(\mathcal{A})$ , where process  $Sup$  collects all alert signals in  $\mathcal{A}$  and converting them in alert signals in  $\mathcal{B}$ , on fresh communication channels.

**Lemma 2.** Let  $M \in \mathcal{M}$  be a CPS with attacker's goal indicators  $\mathcal{I} = \{(i_1, r_1), \dots, (i_k, r_k)\}$ , alert signals  $\mathcal{A} \subseteq \text{out}(M)$ , and detection policy  $\mathcal{P}$ . Let  $Sup$  be a pure-logical process with detection policy  $\mathcal{P}_{Sup}$ , collecting all alert signals in  $\mathcal{A}$  and converting them in alert signals in  $\mathcal{B}$ , on fresh communication channels. Assume that  $\mathcal{P}_{Sup} \sim \mathcal{P}$  for some operator  $\sim \in \{\leq, =, \geq\}$ . Then, for  $\mathcal{P}' = (\mathcal{P} \parallel \mathcal{P}_{Sup}) \setminus \text{chn}(\mathcal{A})$ , the following two properties

$$1) M \xrightarrow{q_1, \dots, q_k} \xrightarrow{\mathcal{P}} \gamma_1$$

$$2) (M \parallel Sup) \setminus \text{chn}(\mathcal{A}) \xrightarrow{q'_1, \dots, q'_k} \xrightarrow{\mathcal{P}'} \gamma_2$$

are equivalent when  $q'_m \sim q_m$ , for  $m \in 1..k$ , and  $\gamma_1$  and  $\gamma_2$  are two distributions such that for any CPS  $M'$  there exists a pure logical process  $Sup'$  for which it holds that  $\gamma_1(M') = \gamma_2((M' \parallel Sup') \setminus \text{chn}(\mathcal{A}))$ .

*Proof.* We prove that 1) implies 2); the proof of the other implication is similar. By Definition 9, the transition in 1) is derived from a trace of the form  $M \xrightarrow{t} \xrightarrow{\text{tick}} \gamma_1$ , with  $t = \alpha_1 \cdots \alpha_n$  and  $\alpha_1, \dots, \alpha_n \in \mathbf{A} \setminus \{\tau, \text{tick}\}$ , such that for all  $m \in 1..k$  we have  $q_m = \min(1, \sum_{j=1, \alpha_j \in \text{out}(M)}^n \mathcal{P}(\alpha_j, i_m))$ . Let us assume  $\text{restr}_{\mathcal{A}}(t) = \overline{c_1}m_1 \cdots \overline{c_h}m_h$  and  $\text{restr}_{(\mathbf{A} \setminus \mathcal{A})}(t) = \beta_1 \cdots \beta_{n-h}$ . As  $\mathcal{P}_{Sup} \sim \mathcal{P}$ , it follows that:

1) there is an execution trace  $Sup \xrightarrow{t'} \xrightarrow{\text{tick}} Sup'$  such that  $t' = c_1m_1 \cdots c_hm_h \cdot \overline{d_1}m'_1 \cdots \overline{d_l}m'_l$ ;

2) for  $m \in 1..k$ , we have  $\min(1, \sum_{j=1}^l \mathcal{P}_{Sup}(\overline{d_j}m'_j, i_m)) \sim \min(1, \sum_{j=1}^h \mathcal{P}(\overline{c_j}m_j, i_m))$ .

Since processes evolve to point distributions, the execution trace  $Sup \xrightarrow{t'} \xrightarrow{\text{tick}} Sup'$  is actually the following:  $Sup \xrightarrow{t'} \xrightarrow{\text{tick}} \overline{Sup'}$ . Thus, from  $M \xrightarrow{t} \xrightarrow{\text{tick}} \gamma_1$  we derive  $(M \parallel Sup) \setminus \text{chn}(\mathcal{A}) \xrightarrow{t''} \xrightarrow{\text{tick}} \gamma_2$ , for a trace  $t''$  of the form  $\beta_1 \cdots \beta_{n-h} \overline{d_1}m'_1 \cdots \overline{d_l}m'_l$ , and  $\gamma_2$  satisfying  $\gamma_1(M') = \gamma_2((M' \parallel Sup') \setminus \text{chn}(\mathcal{A}))$ . Then, by Definition 9 we derive  $(M \parallel Sup) \setminus \text{chn}(\mathcal{A}) \xrightarrow{q'_1, \dots, q'_k} \xrightarrow{\mathcal{P}'} \gamma_2$ , for some  $q'_1, \dots, q'_k$ . In order to complete the proof, it remains to prove that  $q'_m \sim q_m$ , for all  $m \in 1..k$ . The result is obtained by an application of the following algebraic reasoning:

$$\begin{aligned} q'_m &= \min(1, \sum_{j=1}^l \mathcal{P}'(\overline{d_j}m'_j, i_m) + \sum_{\substack{\beta_j \in \text{out}((M \parallel Sup) \setminus \text{chn}(\mathcal{A})) \\ j=1 \dots n-h}} \mathcal{P}'(\beta_j, i_m)) \\ &= \min(1, \sum_{j=1}^l \mathcal{P}_{Sup}(\overline{d_j}m'_j, i_m) + \sum_{\substack{\beta_j \in \text{out}(M) \\ j \in \{1, \dots, n-h\}}} \mathcal{P}(\beta_j, i_m)) \\ &\sim \min(1, \sum_{j=1}^h \mathcal{P}(\overline{c_j}m_j, i_m) + \sum_{\substack{\beta_j \in \text{out}(M) \\ j \in \{1, \dots, n-h\}}} \mathcal{P}(\beta_j, i_m)) \\ &= \min(1, \sum_{\substack{\alpha_j \in \text{out}(M) \\ j=1, \dots, n}} \mathcal{P}(\alpha_j, i_m)) \\ &= q_m \end{aligned}$$

where: step 1 follows from  $t'' = \beta_1 \cdots \beta_{n-h} \overline{d_1}m'_1 \cdots \overline{d_l}m'_l$ ; step 2 follows because  $\mathcal{P}' = (\mathcal{P} \parallel \mathcal{P}_{Sup}) \setminus \text{chn}(\mathcal{A})$ , thus  $\mathcal{P}_{Sup}(\overline{d_j}m'_j) = \mathcal{P}'(\overline{d_j}m'_j)$  and  $\mathcal{P}(\beta_j, i_m) = \mathcal{P}'(\beta_j, i_m)$ , where the last equality holds because no  $\beta_j \in \text{out}(M)$  uses channels in  $\text{chn}(\mathcal{A})$ ; step 3 follows by  $\min(1, \sum_{j=1}^l \mathcal{P}_{Sup}(\overline{d_j}m'_j, i_m)) \sim \min(1, \sum_{j=1}^h \mathcal{P}(\overline{c_j}m_j, i_m))$ ; step 4 derives because  $\text{restr}_{\mathcal{A}}(t) = \overline{c_1}m_1 \cdots \overline{c_h}m_h$  and  $\text{restr}_{(\mathbf{A} \setminus \mathcal{A})}(t) = \beta_1 \cdots \beta_{n-h}$ ; step 5 follows because  $t = \alpha_1 \cdots \alpha_n$ .  $\square$

Now, everything is in place to prove Theorem 2.

*Proof of Theorem 2.* We give the proof for the metric **FP**. The proof for the metric **FN** is similar. We reason by induction on  $n \in \mathbb{N}$ . The base case,  $n = 0$ , is immediate. Let us consider the inductive step,  $n > 0$ . Here, we write  $N$  to denote the CPS  $(M \parallel Sup) \setminus \text{chn}(\mathcal{A})$ , and  $N'$  for CPSs of the form  $N' = (M' \parallel Sup') \setminus \text{chn}(\mathcal{A})$ , where  $M' \in \text{der}(M)$  and  $Sup' \in \text{der}(Sup)$ . The desired result derives from the following algebraic reasoning:

$$\begin{aligned} &\mathbf{FP}_{\mathcal{I}, \mathcal{P}'}^n(N) \\ &= \min \left( 1, \max_{N \xrightarrow{q_1, \dots, q_k} \xrightarrow{\mathcal{P}'} \gamma} \frac{\sum_{j=1}^k r_j \cdot \max(0, q_j - i_j(N)) + (n-1) \cdot \sum_{N' \in \text{supp}(\gamma)} \gamma(N') \cdot \mathbf{FP}_{\mathcal{I}, \mathcal{P}'}^{n-1}(N')}{n} \right) \\ &\sim \min \left( 1, \max_{N \xrightarrow{q_1, \dots, q_k} \xrightarrow{\mathcal{P}'} \gamma} \frac{\sum_{j=1}^k r_j \cdot \max(0, q_j - i_j(M)) + (n-1) \cdot \sum_{N' \in \text{supp}(\gamma)} \gamma(N') \cdot \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^{n-1}(M')}{n} \right) \\ &\sim \min \left( 1, \max_{M \xrightarrow{q_1, \dots, q_k} \xrightarrow{\mathcal{P}} \gamma} \frac{\sum_{j=1}^k r_j \cdot \max(0, q_j - i_j(M)) + (n-1) \cdot \sum_{M' \in \text{supp}(\gamma)} \gamma(M') \cdot \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^{n-1}(M')}{n} \right) \\ &= \mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M) \end{aligned}$$

where: step 1 follows by the definition of  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}'}^n(N)$ ; step 2 follows because  $i_j(M) = i_j(N)$  and by an application of the inductive hypothesis; step 3 follows by Lemma 2; step 4 follows by the definition of  $\mathbf{FP}_{\mathcal{I}, \mathcal{P}}^n(M)$ .  $\square$

## B. Proofs of Section IV

The proofs of Propositions 4, 5, 6 and 7 follow from Theorem 1, Theorem 2 and Proposition 1.