Multi-Dimensional Invariant Detection for Cyber-Physical System Security

A case study of smart meters and smart medical devices

by

Maryam Raiyat Aliabadi

B.Sc., Shahid Beheshti University, 2000 M.Sc., Tehran Islamic Azad University, 2006

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF

MASTER OF APPLIED SCIENCE

in

The Faculty of Graduate Studies

(Electrical & Computer Engineering)

THE UNIVERSITY OF BRITISH COLUMBIA

(Vancouver)

April 2018

© Maryam Raiyat Aliabadi 2018

Abstract

Cyber-Physical Systems (CPSes) are being widely deployed in securitycritical scenarios such as smart homes and medical devices. Unfortunately, the connectedness of these systems and their relative lack of security measures makes them ripe targets for attacks. Specification-based Intrusion Detection Systems (IDS) have been shown to be effective for securing CPSs. Unfortunately, deriving invariants for capturing the specifications of CPS systems is a tedious and error-prone process. Therefore, it is important to dynamically monitor the CPS system to learn its common behaviors and formulate invariants for detecting security attacks. Existing techniques for invariant mining only incorporate data and events, but not time. However, time is central to most CPS systems, and hence incorporating time in addition to data and events, is essential for achieving low false positives and false negatives.

This thesis proposes ARTINALI : A Real-Time-specific Invariant iNference ALgorIthm, which mines dynamic system properties by incorporating time as a first-class property of the system. We build ARTINALI-based Intrusion Detection Systems (IDSes) for two CPSes, namely smart meters and smart medical devices, and measure their efficacy. We find that the ARTINALI-based IDS significantly reduces the ratio of false positives and false negatives by 16 to 48% (average 30.75%) and 89 to 95% (average 93.4%) respectively over other dynamic invariant detection tools. Furthermore, it incurs about 32% performance overhead, which is comparable to other invariant detection techniques.

Lay Summary

Cyber-physical systems (CPSes) constitute the core of Internet of Things (IoT). Recently, they are increasingly deployed in many critical infrastructures. The rapid growth of IoT has led to deployment of CPSes without adequate security. Researchers have demonstrated successful attacks against CPSes used in smart grids, modern cars, unmanned aerial vehicles and smart medical devices. Hence, it is imperative to develop techniques to improve security of these systems. However, CPSes have constraints (such as realtime requirements) that make building Intrusion Detection System (IDS) for them challenging. In this thesis, we propose a technique (ARTINALI) to dynamically monitor the CPS system to learn its common behaviors in three important dimensions including *data*, *event* and *time*, and formulate behavioral rules (aka invariants) for detecting security attacks. Furthermore, we built ARTINALI-based IDSes for two important CPSes, namely *smart meters* and *smart medical devices* and evaluated the efficacy of the IDSes.

Preface

This thesis is the result of work carried out by myself, in collaboration with my advisor, Dr. Karthik Pattabiraman, Dr. Julien Gascon Samson, and Amita Kamath. All chapters are based on work published in ACM SIG-SOFT Symposium on the Foundations of Software Engineering (FSE2017). I was responsible for writing all the papers, designing solutions, conducting the experiments and evaluating the results. My advisor was responsible for editing and writing segments of the manuscripts, providing feedback and guidance during design phases, experiments, and evaluations. Amita contributed primarily to the evaluation of one part of the related work for comparison purposes during her 3-month internship, and Julien helped me with a couple of iterations over my conference paper, and providing useful comments.

"ARTINALI: Dynamic invariant detection for Cyber-Physical System Security." Maryam Raiyat Aliabadi, Amita Ajith Kamath, Julien Gascon-Samson, and Karthik Pattabiraman, In the Proceedings of the ACM SIG-SOFT Symposium on the Foundations of Software Engineering (FSE2017).

Table of Contents

Ał	ostra	ct.		•	•		•	•	•	•	•	•	•	ii
La	y Su	mmar	y	•	•		•	•			•		•	iii
Preface iv														iv
Table of Contents													v	
List of Tables														
List of Figures														
List of Abbreviation														
Ac	Acknowledgments xii													
Dedication														
1	Intr	oducti	on and Overview				•							1
	1.1	Cyber	-Physical Systems		•		•							1
	1.2	Chara	cterizing the CPS Attack Surface .				•						•	2
	1.3	CPS ($Constraints \dots \dots \dots \dots \dots \dots \dots \dots$				•						•	4
	1.4	Intrus	ion Detection Systems				•						•	6
		1.4.1	Signature-based detection $\ . \ . \ .$				•				•		•	6
		1.4.2	Anomaly-based detection				•				•		•	6
		1.4.3	Specification-based techniques				•				•		•	7
	1.5	Overa	rching Goal and Research Questions				•				•		•	9
	1.6	Contri	butions \ldots \ldots \ldots \ldots \ldots \ldots				•				•		•	10

Table of Contents

	1.7	Publications 11
2	Bac	kground and Motivation
	2.1	Related work
	2.2	A Motivating Example 14
	2.3	Summary
3	App	proach
	3.1	Multi-dimensional model
	3.2	Event-data-time Interplay 19
	3.3	ARTINALI Workflow
		3.3.1 Block 1. ARTINALI $D E$ Miner
		3.3.2 Block 2. ARTINALI $E T$ Miner
		3.3.3 Block 3. ARTINALI $D T$ Miner
		3.3.4 Block 4. IDS Prototype
	3.4	Summary
4	Exp	perimental Setup
	4.1	Case Studies 29
		4.1.1 Advanced Metering Infrastructure (AMI) 29
		4.1.2 Smart Artificial Pancreas (SAP)
	4.2	Experimental Procedure
5	Eva	luation Against Targeted Attacks
	5.1	AMI Attacks
		5.1.1 Synchronization tampering (Blocks $A1 - A4$) 34
		5.1.2 Meter spoofing (Blocks $B1 - B5$)
		5.1.3 Message dropping (Blocks $C1 - C5$)
	5.2	Detection of AMI attacks
		5.2.1 Synchronization tampering
		5.2.2 Message dropping 36
		5.2.3 Meter spoofing $\ldots \ldots \ldots \ldots \ldots \ldots \ldots 37$
	5.3	SAP Attacks
		5.3.1 CGM spoofing attack (Blocks $A1 - A4$)

		5.3.2 Basal tamperin	g (Blocks $B1 - B5$)	38
	5.4	Detection of example a	attacks in SAP	39
		5.4.1 CGM spoofing	attack	39
		5.4.2 Basal tamperin	g attack	39
		5.4.3 Summary	~	40
6	Eva	luation Against Arbi	trary Attacks	41
	6.1	Arbitrary Attack Mode	el	41
	6.2	Evaluation Metrics		43
	6.3	Research Questions (R	Qs)	45
		6.3.1 RQ1. F-Score	· · · · · · · · · · · · · · · · · · ·	45
		6.3.2 RQ2. False Neg	gatives	51
		6.3.3 RQ3. False Pos	itives	53
		6.3.4 RQ4. Memory	Overhead	54
		6.3.5 RQ5. Performa	nce Overhead	55
	6.4	Summary		56
7	Dise	cussion		58
	7.1	Threats to Validity .		58
	7.2	Generalizability of AR	TINALI	59
8	Con	clusions and Future	Work	60
	8.1	Summary		60
	8.2	Future work		61
		8.2.1 Extending the	Generalizability of ARTINALI	62
		8.2.2 Optimizing AF	RTINALI for the Scalability	62
		8.2.3 Extending AR	TINALI for Attack Diagnosis	63
Bi	bliog	raphy		64

List of Tables

1.1 The main deficiencies of current IDS techniques in address					
	CPS constraints	8			
3.1	$\mathbf{E} \mathbf{T} \text{ and } \mathbf{D} \mathbf{T} \text{ Invariant Types } \ldots \ldots \ldots \ldots \ldots \ldots \ldots$	25			
4.1	Types and number of inferred invariants for SEGMeter across				
	tools	33			
4.2	Type and number of inferred invariants for OpenAPS across				
	tools	33			
5.1	ARTINALI invariants to detect the example attacks in AMI.	36			
5.2	ARTINALI invariants to detect example attacks in OpenAPS.	38			
6.1	Optimal training set size for maximum $F-Score(0.5)$ for SEG-				
	Meter across tools, and the corresponding FP and FN ratios.	47			
6.2	Optimal training set size for maximum F-Score(2) for Ope-				
	nAPS across tools, and the corresponding FP and FN ratios.	47			
6.3	Memory and performance overhead of IDS, seeded by ARTI-				
	NALI and the other tools, running on SEGMeter	55			

List of Figures

1.1	CPS attack surface	3
2.1	Scope of dynamic invariant detection techniques	14
2.2	A snippet of Serial-Talker code for the SEGMeter	15
3.1	Work flow of ARTINALI	22
3.2	Sample DElog and respective $D E$ invariants with 100% con-	
	fidence	24
3.3	Sample TElog and DurationLOG files and respective $E T$ in-	
	variants	26
3.4	Representation of ARTINALI $D T$ invariants $\ldots \ldots \ldots$	27
4.1	High level block diagram of (a) Smart meter and (b) Smart	
	Artificial Pancreas.	31
4.2	Overall experimental process of running the IDS	32
5.1	Attack tree for AMI	35
5.2	Attack tree for SAP	38
6.1	Fault injection $impact(\%)$: (a) data mutation in SEGMeter,	
	(b) branch flip in SEGMeter, (c) artificial delay in SEGMeter,	
	(d) data mutation in OpenAPS, (e) branch flip in OpenAPS,	
	(f) artificial delay in OpenAPS	43
6.2	FN, FP and F-Score variations based on number of training	
	traces for SEGMeter's IDS seeded by Daikon invariants. $\ $	47
6.3	FN, FP and F-Score variations based on number of training	
	traces for SEGM eter's IDS seeded by Texada invariants. $\ $	48

List of Figures

6.4	FN, FP and F-Score variations based on number of training	
	traces for SEGM eter's IDS seeded by Perfume invariants. $\ . \ .$	48
6.5	FN, FP and F-Score variations based on number of training	
	traces for SEGMeter's IDS seeded by ARTINALI invariants.	49
6.6	FN, FP and F-Score variations based on number of training	
	traces for OpenAPS's IDS seeded by Daikon invariants. $\ . \ .$	49
6.7	FN, FP and F-Score variations based on number of training	
	traces for OpenAPS's IDS seeded by Texada invariants	50
6.8	FN, FP and F-Score variations based on number of training	
	traces for OpenAPS's IDS seeded by Perfume invariants. $\ . \ .$	50
6.9	FN, FP and F-Score variations based on number of training	
	traces for OpenAPS's IDS seeded by ARTINALI invariants.	51
6.10	$\mathrm{FN}(\%)$ of IDS for SEGMeter for different attack types across	
	the tools. Error bars are shown for the 95% confidence inter-	
	val	52
6.11	$\mathrm{FN}(\%)$ of IDS for OpenAPS for different attack types across	
	the tools. Error bars are shown for the 95% confidence inter-	
	val	52

List of Abbreviation

- \bullet $\mathbf{AMI}:\mathbf{A} \mathrm{dvanced}$ $\mathbf{M} \mathrm{etering}$ $\mathbf{I} \mathrm{n} \mathrm{frastructure}$
- $\bullet \ \mathbf{BG}: \mathbf{B} \mathrm{lood} \ \mathbf{G} \mathrm{lucose}$
- \mathbf{FNR} : False Negative Ratio
- $\bullet~\mathbf{FPR}:\mathbf{False}~\mathbf{Positive}~\mathbf{R}atio$
- $\bullet~\mathbf{IDS}:\mathbf{Intrusion}~\mathbf{D}\mathbf{e}\mathbf{tection}~\mathbf{S}\mathbf{y}\mathbf{s}\mathbf{tem}$
- IoT : Internet of Things
- PCD : Power Consumption Data
- $\bullet~\mathbf{SAP}:\mathbf{S}\mathrm{mart}~\mathbf{A}\mathrm{rtificial}~\mathbf{P}\mathrm{ancreas}$
- **SDC** : Silent Data Corruption
- $\bullet~\mathbf{UAV}:\mathbf{Unmanned}~\mathbf{Aerial}~\mathbf{V}ehicle$

Acknowledgments

Firstly, I would like to express my sincere gratitude to my adviser, *Prof. Karthik Pattabiraman*. If it was not for your continuous support, careful guidance, and great patience, I would not have been able to complete this thesis today.

Besides my adviser, I would like to thank my thesis committee and examiners, including *Prof. Ivan Beschastnikh*, and *Prof. Sathish Gopalakrishnan*. You provided me with valuable feedback, and helped me improve my thesis.

I would also like to thank my colleagues in CSRG for their help and insightful feedback on my work. I particularity would like to thank all my colleagues at Dependable Systems Lab for the stimulating discussions, and all the joy we had working together. The memories will stay with me forever.

My lovely *ARTIN* and my dearest *ALI*, since starting my studies at UBC and living in Vancouver, you have supported me unconditionally, and beyond imagination. I am very lucky to have you in my life, and I will be forever grateful to you.

My loving *parents* and *parents-in-law*, I would like to thank you for giving me your wonderful support through ups and downs. You shaped me into who I am today, and helped me grow both as a researcher, and as a person.

Last but not the least, I would like to thank *God* for giving me the chance to pursue this wonderful profession, and for guiding me throughout the entire process.

Dedication

To my son, ARTIN and my husband, ALI

Chapter 1

Introduction and Overview

1.1 Cyber-Physical Systems

Cyber-physical systems (CPSes) have been investigated as a key area of research since they constitute the core of the Internet of Things (IoT). Recently, they are increasingly deployed in many critical infrastructures. For example, they are widely used in modern automotives to control different parts of a car [10]. Millions of people throughout the world have implanted medical smart devices such as pacemakers to help their hearts beat at a regular rhythm [29]. Moreover, recent studies predict that global spending on Unmanned Aerial Vehicles (UAVs), will exceed \$94 billion over the next ten years [53] and the number of Advanced Metering Infrastructures (AMIs) in smart grids will exceed 800 million by 2020 [49].

The rapid growth of IoT has led to deployment of CPSes without adequate security. These systems carry out critical tasks and hence, are potential targets for cyber attacks. For instance, hackers are able to remotely kill the engine of the smart car, and take control of the brakes [26], or cause a drone to crash or go off course[21, 24, 43, 47]. The U.S. Department of Homeland Security investigated 24 cases of suspected cyber security attacks in smart medical devices in 2014 alone[29, 33], and this number is expected to increase. In addition, researchers have successfully discovered and demonstrated a variety of attacks in smart meters for energy fraud purposes [45]. To make these systems secure, security mechanisms such as intrusion detection systems (IDS) are required.

1.2 Characterizing the CPS Attack Surface

A CPS consists of a cyber unit (i.e., embedded system), and a physical unit connected by a communication channel. It also includes a control loop which involves interactions between the cyber and physical domains. In Figure 1.1, we characterized the attack surface of a CPS, and extended the attacker entry points presented by previous work [9, 28, 34].

According to Figure 1.1, there are four likely entry points (marked as A, B, C, D) for attackers to penetrate the system.

Type A attacks can be used to hide the real state of the physical process in order to trick the controller program to make a wrong decision about the next state of physical process, and to delay the detection of the attacks before the actual damage to the system (like what happened in *Stuxnet* [11]).

Type C attacks can disrupt the system by directly compromising the control commands that are issued by controller. These attacks can jam the communication channel, or compromise the routing protocol. One example of this type is *Basal tampering attack*, which transmits the malicious commands to the smart insulin pump after crafting a DoS attack by jamming the communication between controller and insulin pump [33, 41].

Type B include attacks that exploit the vulnerabilities in physical components. For example, [47] developed a way for attacking drones equipped with vulnerable MEMS gyroscopes using intentional sound noise causing drones to loose control and crash.

Type D attacks aim to exploit the vulnerabilities in cyber unit to take over the dynamics of physical process. For instance, recent work [48] investigated an attack to smart facial recognition systems caused by exploiting a mis-classification bug (CVE-2016-1516) in the controller algorithm.

In this study, we focus on attacks that impact the physical processes without exploiting vulnerabilities in the physical domain (type A, C and Dattacks). The reason is that many physical systems deal with legacy technology and legacy components, that are not controlled by software, and their behavior is essentially governed by the physical law. Hence, it is challenging to identify the environment variables, map the environment changes on the



Figure 1.1: CPS attack surface

environment variables, and incorporate the laws of physics to define acceptable behavior upon environmental changes [37]. Therefore, the behavior of physical processes can be defined more precisely by specifications that are specified by system developer or experts' knowledge.

Adversarial Model : We assume the adversarial goal is to compromise the functionality of CPS. This means that she either prevents a vital functionality of the CPS from being executed (e.g., power consumption data not being sent to the utility server in smart meter, or regular heartbeat not being provided by a pacemaker), or functionalities being executed improperly (e.g., insulin injection is resumed when it must be stopped in insulin pump, or brake leads to acceleration in a smart car). The first group of attacks targets the availability properties, while the second group targets the integrity properties of the CPS. We assume that the adversary is capable to inject, drop or modify messages in the communication channels on attack entry points A, C and D defined in previous section to change i) the legal values of data variables in the code , ii) the normal execution path in the control flow graph of the code, or iii) the normal timing behavior of a particular function of the code, and hence impact the CPS functionality. We do not consider denial of service (DoS) attacks or those threats that compromise the privacy/confidentiality properties of the CPS, as these attacks are generally addressed by other techniques including network security measures and cryptographic methods.

1.3 CPS Constraints

In the following, we discuss the key characteristics and constraints that security mechanisms must satisfy to be applicable to CPSes.

- C1. Real-time constraints: CPSes typically interact with their environment in a real-time fashion. From a security point of view, taking real-time requirements into account is vital for two reasons: First, CPSes are decision making agents that need to make decisions in real time, and to address the continuous operation capabilities; therefore, real-time availability is a necessity. This characteristic leads to a stricter operational environment [8], in which any security solution that modifies the controller program and changes its real-time behavior is not acceptable for these systems. In other words, the performance overhead imposed by security mechanism has to be small enough to satisfy the CPS real-time constraints. Secondly, in a real-time system, the operational correctness depends on both logical correctness, and correct timing behavior [54]. This implies that the logical assertions are not enough for verifying the correctness of these systems, and hence incorporating time in their system model is essential to enable the security mechanism to check the system's operations [15]. In other words, reflecting the real-time constraints into the CPS model is required to have a more predictive security mechanism.
- C2. Resource constraints : As a single thing in the "Internet of Things", a CPS performs a single task on a single platform with limited CPU, memory and computational power. It implies that these systems need a security solution that satisfies their resource constraints. For instance, an important component of a security mechanism is the

model that represents the correct behavior of the CPS. This model may occupy a large space in memory. However, CPSes have limited memory capacity making existing IDSes inapplicable due to memory overheads.

- C3. Unknown vulnerabilities : CPSes are new systems that may have unknown vulnerabilities, and hence, they are inevitably exposed to zero-day attacks. On the other hand, the physical system is more susceptible to the vulnerabilities of the cyber system as the interaction between the physical system and the cyber system increases. Thus, CPSes require security solutions that rely on a system model rather than a dictionary of attacks, to monitor both known and unknown attacks.
- C4. Security-critical constraints : As many CPSes are deployed in security-critical applications such as smart medical devices, they need security mechanisms with minimum false negative rates. This is important because any undetected attack in such life-critical systems may have catastrophic consequences for the patient.
- C5. Large-scale deployment : CPSes are often deployed on a large scale (e.g., smart meters), and hence need security mechanisms with minimal false positives as false positives can aggregate within the system and consume significant resources. Moreover, these systems are deployed in mission-critical applications where shutting them down on a false alarm is not a viable option.
- C6. No-human-in-the-loop : Most CPSes are autonomic systems, which work without a human-in-the-loop, and need to provide non-interrupted service. Thus, unlike the general computer systems, CPSs cannot be interrupted frequently for upgrading and/or patching. This implies that there is a substantial need for deployment of highly sensitive security mechanisms that provide automated response against security attacks. Particularly, this is essential for those CPSes that

are deployed in the security-critical applications such as pacemakers, where any service interruption may be fatal for the pacemakerimplanted patient, or those that are deployed on a large scale including smart meters, where denying the service due to an attack may shut down the entire smart grid [50].

1.4 Intrusion Detection Systems

Intrusion Detection Systems (IDSes) are the most popular security mechanisms that have been widely used to monitor computer systems and detect security attacks. Typical IDSes fall into one of three categories: *Signaturebased*, *Anomaly-based*, and *Specification-based* [20]. In the following, we first introduce the existing IDSes used for computerized systems, then we explain why they are not sufficient to secure CPSes (as illustrated in Table 1.1).

1.4.1 Signature-based detection

Signature-based detection techniques work by monitoring system behavior and comparing the behavior against a database of signatures or attributes from known malicious threats. The benefit of these techniques is that they do not need a model of the system they are monitoring, and their limitation is that they cannot detect zero-day attacks [38]. The latter is especially important for CPSes as they are often difficult to patch or upgrade in the field. Therefore, signature-based IDSes are not a good match for CPSes as they do not address constraints C3 and C6. In contrast, both anomaly-based and specification-based techniques use a behavioral model of the system to compare with suspicious behaviors, and can detect both known and unknown attacks.

1.4.2 Anomaly-based detection

Anomaly-based techniques create a baseline profile of the legitimate system behavior based on statistical methods by observing its operations at runtime, enabling them to distinguish the incoming system activity as normal or anomalous. Unfortunately, these techniques incur considerable overhead to profile the system at runtime [20], and also suffer from high rates of false-positives, both of which deter their use in CPSes, which are often resource constrained, and operate autonomously for long periods of time. Accordingly, they are not viable for CPSes with respect to constraints C2, C4 and C5.

1.4.3 Specification-based techniques

Specification-based techniques build a behavioral model for system by defining a set of rules (known as *invariants*) that lead to a decision regarding whether a given pattern of activity is suspicious. Invariants are defined as statements that describe the relationship among states of a program [4]. For example, they can be used to state the relationship between two variables that always hold true, over the entire run of a program. Invariants can be either discovered by *static analysis* or *dynamic analysis* of the program, or be provided by the *system developer*.

Static analysis-based techniques :Static analysis-based techniques are based on source code and the specifications defined by the developer. Hence they rely upon apriori knowledge of the systems' specifications to detect attacks [7]. Static analysis-based techniques are inherently conservative with low false positives, as they only generate invariants that are rigorously provable. These techniques analyze the program without executing the program, and hence, they can not provide adequate information about the real-time behavior of the system in its operational environment. Thus, static analysis-based techniques are not able to provide a comprehensive model of real-time behavior and timing properties of system, which in turn, leads to high false negatives [20]. Moreover, as these techniques need to access the source code, they cannot statically verify many interesting properties of a program because of unavailable code (e.g., calls to compiled external libraries). As a result, the invariants set that is generated by static analysis techniques needs to be complemented by developer specifications to address certain external conditions [4]. However, prior work [14, 39] have found that there is of-

Table 1.1: The main deficiencies of current IDS techniques in addressing CPS constraints

IDS Techniques	C1	C2	C3	C4	C5	C6
Signature-based			\checkmark			\checkmark
Anomaly-based		\checkmark		\checkmark	\checkmark	
Specification-based (static analysis)	\checkmark			\checkmark		
Invariant-based (dynamic analysis)				\checkmark	\checkmark	

ten inconsistency between what developers describe their system does, and what the system does in practice [14, 39]. These drawbacks make the static analysis-based techniques insufficient for CPSes with respect to constraints C1 and C4.

Dynamic analysis-based techniques :Dynamic analysis-based techniques provide an alternative way to understand the system by observing the runtime behavior. They log the key points of the program to peek into the actual program behavior at run-time[40], and infer a set of likely invariants. Unlike static analysis-based techniques, these techniques do not need the source code as they rely on data received from runs of program. Although they instrument the code to collect data for analysis they do not need the code for analysis step by itself. Using a sufficiently complete set of test cases, dynamic analysis-based techniques are potentially able to infer invariants that cannot be mathematically proved (e.g, time duration between two specific events of the program), so there is no way for static analysis-based techniques to find them. For these reasons, to address CPS constraints discussed in the previous section, we selected dynamic analysis for mining likely invariants in CPS systems.

There has been a significant amount of work on using dynamic analysis to find likely invariants for program understanding, formal verification, debugging and intrusion detection [12, 17, 19, 23, 27, 30, 31, 35, 40, 42, 55, 58]. These systems mine execution traces of the system for deriving invariants on the data values of the system, the events or both. However, we find that many of these systems incur significant false-positives and/or false-negatives, when used in the context of an IDS, which makes them challenging to deploy for CPSes that are used in security critical infrastructures (C4) and/or on a large scale (C5).

1.5 Overarching Goal and Research Questions

CPSes are targets of many security attacks. Reconfiguration and/or recovery of such systems from attacks requires time, money, and human effort. On the other hand, as a result of what we discussed in previous section, the existing IDSes for computer systems do not address the limitations and constraints of CPSes. Therefore, the overarching goal of this thesis is to understand the behavior of CPS systems in the absence of attacks, and use this understanding to develop an IDS technique to improve the security of CPS systems in an automated manner.

To achieve the overarching goal, we are interested in answering the following research questions:

RQ1. How can we build a specification-based IDS for CPS systems with respect to their constraints?

- RQ1.1. How can we infer a multi-dimensional behavioral model for CPS systems?
- RQ1.2. How can we leverage the multidimensional CPS model to build an IDS that meets CPS constraints?

The most important component of a specification-based IDS for CPS systems is the CPS model. A CPS model should include a set of specifications (invariants) that precisely describe the correct behavior of the system. We addressed these research questions by proposing a technique to analyze the normal behavior of the CPS along three dimensions: data, event, and time, to generate a richer model for the system, and hence a more sensitive IDS for the system.

1.6 Contributions

This thesis introduces ARTINALI (A Real-Time-specific Invariant iNference ALgorIthm) for mining likely invariants through dynamic analysis in CPS systems, for specification-based IDSes. The main innovation in ARTINALI is that it incorporates time as a first-class notion in the mined invariants, in addition to the traditional data and event invariants. This is important for two reasons. First, most CPSes have real-time constraints, and hence their operational correctness depends on both logical correctness, and correct timing behavior [25, 54]. Hence, incorporating time is essential for detecting many common security attacks in these systems. Secondly, CPSes have predictable timing behaviors to a first order of approximation, and hence leveraging this predictability leads to higher accuracy (i.e., lower false-positives and negatives). However, incorporating time in dynamic invariant detection techniques increases the complexity of the learning due to the much larger state space that needs to be covered. To alleviate this issue, we break up the problem of learning invariants along the three dimensions into problems of learning invariants along two dimensions, namely data-events and eventstime, and then combine them into data-events-time invariants. To the best of our knowledge, ARTINALI is the first dynamic invariant detection system that mines invariants along the three dimensions of data. event. and time, and uses the mined invariants for intrusion detection.

Our contributions are:

- Designed ARTINALI, an algorithm that generates a multi-dimensional model for CPSes by mining invariants along the data, event and time dimensions (Chapter 3).
- Built an ARTINALI-based IDS prototype, and used it in the context of two CPS systems, namely i) advanced metering infrastructures, ii) and smart artificial pancreas (Chapter 4).
- Evaluated our ARTINALI-based IDS prototype on the two CPS systems, in comparison with several existing state of the art dynamic invariant detection techniques. We crafted 6 targeted attacks against

two systems, and observed that ARTINALI-based IDS is able to detect all of them while the other techniques could not detect even a single attack (Chapter 5).

• Found that the ARTINALI-based IDS exhibits significantly lower falsenegatives and false-positives for arbitrary attacks emulated by fault injection, compared to the other techniques. Furthermore, it incurs about 32% performance overhead, which is comparable to other invariant detection techniques (Chapter 6).

1.7 Publications

The work in this thesis has been published in the following paper:

• Maryam Raiyat Aliabadi, and Karthik Pattabiraman. "ARTINALI: Dynamic invariant detection for Cyber-Physical System Security." ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE2017).

Chapter 2

Background and Motivation

Specification based techniques based on static analysis and formal specifications (generated by developer) have been proposed as a good fit for CPS security [7, 37]. However, as we discussed in chapter 1-Section 1.4.3, static analysis-based techniques have deficiencies that make them insufficient for addressing the CPS constraints (e.g., real-time constraints). On the other hand, developers rarely write down specifications of their systems. As a result, many specification mining techniques based on dynamic analysis have been appeared in previous work [17, 19, 23, 27, 30, 31, 35, 40, 42, 55, 58]. Mined specifications cannot replace formal specifications created by an expert since an invariant mined from program traces could be a false positive [3]. However, as many CPS programs lack formal specifications, mined specifications are necessary. As a result, dynamic analysis is substantially important to mine invariants that can not be inferred through static analysis or are not specified by system developers. For these reasons, we selected to employ dynamic analysis to formulate CPS behavior with respect to their constraints for specification-based IDSes.

In this Chapter, we first survey related work in the area of dynamic analysis-based techniques and explain how our proposed technique differs from them. We then present a motivating example from smart meters to illustrate why we need new types of invariants like the ones generated by our technique to bridge the gap.

2.1 Related work

There has been a significant amount of work on using dynamic analysis to model the behavior of software systems for program understanding, formal verification, debugging and intrusion detection [17, 19, 23, 27, 30, 31, 35, 40, 42, 55, 58]. These techniques can be categorized into four classes, based on the models that they generate: i) data invariants, ii) event relationships, iii) data and event relationships, and iv) time dependencies of events. Figure 2.1 shows the main dynamic analysis-based techniques, and where they fall along the data, event and time axes.

Daikon was the first dynamic analysis-based technique to derive (likely) invariants about data value relations [17], and falls into the first class of techniques. Daikon can be placed on the *data* axis as it produces a model for data constraints without taking into account the events or timing of the system. DySy [13], which uses symbolic execution to derive invariants, is another example of this class.

The second class captures the sequence of events within a progam's execution paths by inferring finite state machines from a set of traces. Relevant examples include Perracotta [58] and Texada [31], both of which derive temporal logic propositions, and capture sequences of events by tracking dynamic traces. These tools fall along the *event* axis since they only capture the constraints on event relations independent of data or timing information.

The third class of techniques generate integration models that capture the relationship between data and events. For example, the GK-Tail algorithm merges temporal specifications and data invariants into Extended Finite State Machine models [35]. It represents sequences of method invocations that are annotated with data, and is hence limited to classifying data invariants that arise among method calls. Quarry finds data invariants at each program point, and then finds temporal relationships between the invariants [30]. Neither technique considers timing information, however.

The fourth class consists of a single technique, Perfume, which is a specification mining tool designed for modelling system properties based on resource (time and storage) consumption [40]. It generates an integration model of event relations and their time constraints. Although Perfume considers time as a part of model, it does not consider the relationship between data and time.

Overall, none of the current techniques consider the interplay among

time, events and data in formulating invariants, which we believe is an essential characteristic of CPS systems.



Figure 2.1: Scope of dynamic invariant detection techniques

2.2 A Motivating Example

We consider an example of a smart electric meter to illustrate why the existing dynamic analysis-based techniques are often insufficient for capturing the key properties of a CPS system. We also use this as a motivating example to illustrate ARTINALI later.

We use an open-source smart meter called SEGMeter as one of the testbeds for our implementation and our evaluations (see Chapter 4-Section 4.1.1). SEGMeter is composed of two main components: the meter component and a controller. The meter component is in charge of measuring and collecting power consumption data coming through its serial ports, and storing them in memory. The controller acts as the communication bridge between the meter board and the server, and is in charge of passing server commands to the meter board, as well as transmitting power consumption data to the server at specific time intervals. The Serial-Talker() function in the controller program of the smart meter is in charge of receiving power consumption data (at specific time intervals) and buffering them for billing calculation purposes. The Serial-Taker code is shown in Figure 2.2 (in the Lua language).

```
function serial_talker()
1
    . . . . .
       seg-data = get_data_timer()
6
       if (seg-data == true ) then
         command = "(all_nodes (start_data))"
         serial client:send(command .. ";\n")
         read(message)
         . . . .
       end
        if (seg-data == false ) then
24
         stream, status, partial = serial client:receive(16768)
         f:write("node_name = ", tostring(partial), "\n")
28
         . . . .
       end
   . . . .
56 end
```

Figure 2.2: A snippet of Serial-Talker code for the SEGMeter

Serial-Talker() has a Boolean argument seg-data, that takes values true or false in pre-determined time intervals. The sequence of events that are invoked in this function varies based on the value passed in argument (seg-data). If true is passed (line 6 in Figure 2.2), then the program emits the event send, followed by read. Alternatively, if false is passed to the function (line 21 in Figure 2.2), then the program emits events receive and write respectively.

We examined the invariants inferred by the different dynamic analysisbased tools for this example. Daikon infers the values of variable seg-data within Serial-Talker() during normal execution as the set $\{true, false\}$, namely seg-data: [true,false]. A typical temporal specification miner such as Texada identifies the legal sequences of events, e.g., $G(\text{send} \rightarrow$ XF read), which means that upon event send happening, it is always followed by event read. The invariant inferred by Perfume (send \rightarrow receive, 0.1, 1.2) complement the temporal invariant by adding time boundaries between events, i.e., send is followed by read within a time interval of 0.1 to 1.2 ms.

Assume that an adversary's goal would be to perform energy fraud and lower their energy bills. One possible attack would be for the attacker to tamper with the synchronization between the **send** and **receive** modes in the smart meter. As a result, a part of the energy usage would not be written to the memory buffer which is used for future energy usage calculations and billing. For instance, should the value **false** be passed to the function instead of **true**, then it would lead to the execution of **receive** and **write** instead of **send** and **read**; hence the billing information would be incorrect.

None of the above techniques can detect the attack as the incorrect occurrence of sequences are triggered by legal values of seg-data occurring at the wrong time (e.g., seg-data (T1) = false). More specifically, Daikon would notice a valid value for seg-data, Texada would notice a normal sequence of receive and write events, and Perfume would also observe valid time intervals between events receive and write within the executed path. Thus, none of them would detect the attack. Even if all three models are used jointly, they would still not detect the intrusion, as the different models either capture the legal data values, or the legal sequence of events with their time difference, but not the interplay among them. This interplay is essential for detecting the attack.

2.3 Summary

In this chapter, we surveyed the dynamic analysis-based techniques that model the behavior of software systems. We categorized these techniques into four classes, based on the models that they generate, and illustrated where they fall along the data, event and time dimensions. Overall, none of these techniques consider the relationship between data and time, as well as interplay among time, events and data in formulating invariants, which we believe are essential characteristics of CPS systems. We also brought up a concrete attack example against a smart electric meter, and showed

2.3. Summary

that none of the existing techniques would detect the attack. Even if all the models generated by different techniques are used jointly, they would still not detect the attack, as the different models either capture the legal data values, or the legal sequence of events with their time difference, but not the interplay among them. This interplay is essential for detecting the attack, and is the main gap in existing dynamic analysis-based techniques.

Chapter 3

Approach

In this chapter, we introduce the security model that ARTINALI uses, and we explain its design. We first define our multi-dimensional model and the different classes of invariants. Next, we explain how to relate different dimensions to generate real-time data invariants. Finally, we present the ARTINALI workflow and algorithm.

3.1 Multi-dimensional model

We model a CPS in three dimensions, as follows:

Data refers to data values assigned to the variables of a program. It includes neither the timing of processes, nor the sequence and concurrency of processes.

Event refers to an action that a system takes to respond to an external stimulus.

Time refers to real-time constraints, and includes both the constraints on physical timing of various operations, and those where the system must guarantee response within a specified time frame.

We model the security policy of a CPS by inferring the set of *invariants* to be preserved during run time. An *invariant*, or interchangeably a *property*, is a logical condition that holds true at a particular set of program points. Like in prior work [17, 31, 58], we use the term *invariants* as a short-cut for *likely invariants*, which are the properties that we *observe* to be true across a set of dynamic execution traces. Corresponding to the dimensions defined above, we define six major classes of invariants that form the basis of the CPS model.

- *Data Invariant* captures the expected range of values of selected data variables during normal execution of program.
- *Event Invariant* captures common patterns in the system's events such as the order of the events' occurrence.
- *Time invariant* captures the normal time boundaries (such as duration or frequency) of an event.
- Data per Event(D|E) invariant captures the temporal relationship between data and events. It allows the IDS to check the validity of data invariants based upon events.
- Event per Time (E|T) invariant captures the constraints over event and time. It represents the boundaries of transition time from one event to another in an event sequence.
- Data per Time (D|T) invariant captures the relational constraints of time and data invariants, or the data invariant as a function of time.

3.2 Event-data-time Interplay

In a CPS, an event is defined as an instance of an action that leads to a change of condition [51] (e.g., message send/receive, sensor data reading, or activating insulin injection). Events have three key features. First, they reflect interactions between system components and observations rather than internal state. The second feature is the notion that events are separated in space and time [15, 16, 51]. Thirdly, the locations in the code where events are triggered are usually *system calls* that are accessible by attackers. From a security perspective, events are important as they play the role of an input channel for malicious communication with the CPS. For instance, those points in which a new measurement is read from sensors, or actuation commands are sent to physical components, are more vulnerable to *spoofing attacks* [18].

Finding a direct relationship between time and data is challenging from both the *learning* and *detection* perspectives. Since time is a continuous

3.2. Event-data-time Interplay

phenomenon, we cannot define a sharp time for transitions in data values or changing states of the system; instead, a distribution of time values has to be learned. As execution time variations might be caused by differences in input sets or different execution flows, rather than malicious activities, the invariant inference technique should learn the normal time variations of the system. On the other hand, the IDS has to distinguish legitimate time variations from any time deviation that indicates an intrusion. To overcome these challenges, we leverage the event-based nature of a CPS, in which every event takes place in an unique time-frame. We discretize the *time* by the *events*, and use these for learning invariants. After discretizing the time by events, we first examine the relationship between data and event dimensions to produce invariants that integrate event information with constraints on data values (D|E invariants). Secondly, we discover the relational constraints over time and event dimensions to calculate the physical time boundaries of events, either independently (time invariants), or in relation to each other (E|T invariants). Finally, we combine the result of the previous steps to infer D|T invariants.

In the following discussion, we illustrate how we infer the D|T invariants given the conditional probability of having *data* D given *event* E invariant (P(D|E)), and given the conditional probability of having *event* E given *time* T invariant (P(E|T)). It should be noted that we have represented different classes of our *likely* invariants in the form of probability functions to conceptually describe the process of how we infer real-time data invariants (i.e., D|T invariants). We exploited the conditional independence of time Tand data D upon a given event E_j to derive equation 3.9. However, ARTI-NALI does not calculate the probability of correctness for every invariant; instead, it generates the invariants that hold true with a high probability based on the above analogy.

Considering data D, event E and time T as random variables, equation 3.1 expresses the joint probability distribution of variables D, E and T. We rewrite it to obtain equation 3.2. From these two equations, we then derive

equation 3.3, which expresses the probability of having D and E, given T.

$$P(D, E, T) = P(D, E|T) \cdot P(T)$$
(3.1)

$$P(D, E, T) = P(D|E, T) \cdot P(E|T) \cdot P(T)$$
(3.2)

$$P(D, E|T) = P(D|E, T) \cdot P(E|T)$$
(3.3)

Using the marginal probability mass function of D shown in equation 3.4, we formalize P(D|T) (the probability of having D given T) in equation 3.5 as the sum of the probabilities of data D and event E_j given time T for all events E_j , which can then be rewritten as equation 3.6 (using equation 3.3).

$$P(D) = \sum P(D, E_j), \forall E_j$$
(3.4)

$$P(D|T) = \sum P(D, E_j|T), \forall E_j$$
(3.5)

$$P(D|T) = \sum P(D|E_j, T) \cdot P(E_j|T), \forall E_j$$
(3.6)

For example, assuming that at time T, event E_j occurs; and that upon E_j occurring, then variable D gets assigned a specific value. This implies that T is the cause of E_j , and that D is the effect of E_j . Thus, variable D is conditionally independent of time variable T given event E_j . Consequently, D and T are conditionally independent, and $P(D|E_j,T) = P(D|E_j)$. Hence, we can simplify the formulation of P(D|T) as follows :

$$P(D|T) = \sum P(D|E_j) \cdot P(E_j|T), \forall E_j$$
(3.7)

According to the event-based semantics of CPS, any given event takes place in a unique time frame. This implies that two or more events cannot take place at the same time T; i.e., $P(E_j|T) > 0 \Rightarrow P(E_i|T) = 0, \forall E_i \neq E_j$. Given this assumption, we first rewrite equation 3.7 to obtain equation 3.8. Then, we simplify it to obtain equation 3.9, which captures the relationship between data D and time T by exploiting the relational constraints of both data and time over the same event E_j which takes place at time T.

$$P(D|T) = P(D|E = E_j) \cdot P(E = E_j|T) + \sum P(D|E_i) \cdot P(E_i|T), \forall E_i \neq E_j$$
(3.8)

$$P(D|T) = P(D|E = E_j) \cdot P(E = E_j|T)$$
(3.9)

In other words, for a given event E_j , a D|T invariant holds true (i.e., happens with a high probability) if and only if both the corresponding D|E invariant and E|T invariant hold true.

3.3 ARTINALI Workflow

ARTINALI is a *dynamic analysis-based* technique that generates models of dynamic system behavior, and proposes a multi-dimensional model based on the design concepts introduced in the previous section. Figure 3.1 shows the key blocks of ARTINALI's workflow.



Figure 3.1: Work flow of ARTINALI

ARTINALI technique works at *event granularity* to mine invariants. In order to generate logs for mining invariants, we manually instrument events

and their associated data variables¹. As we generate the CPS model for attack detection, we capture all system calls (which are reachable by attackers) as events. However, in ARTINALI, events are user-defined. Accordingly, user can optionally customize the level of granularity by choosing another type of events, or prune the space of events by specifying only the *important* system calls based on the system's requirements. We instrument the events' program locations by inserting calls to the *ARTINALI API functions* that we developed for collecting logs, before and after the event. During the runtime, these functions collect *data* and *time* information associated with the instrumented events in separate log files (*DElog* and *TElog*). The logged information is used as the basis for mining invariants.

3.3.1 Block 1. ARTINALI D|E Miner

The ARTINALI D|E Miner learns invariants about the variable values, and how these values relate to a particular event in the system using a threestep process. First, the D|E Miner takes the logged information, and groups them within each trace into distinct classes labeled with the events. It then merges classes across the training traces. Second, within each class, using the Frequent Item Set mining algorithm [22], it merges the data variables while calculating the level of the *confidence* and *support* for every variable. As in prior work [17], *Support* is the fraction of traces in which the variable x within class E_j is seen, and *confidence* is the fraction of supported classes, where variable x is assigned to the same value(s).

Finally, the D|E Miner infers the data invariants associated with each class (event). D|E invariants are multi-propositional data invariants as they all hold true within the same observed event (at the same time). The D|E invariants are stated in the form of $(E_i : d_1 = [], d_2 = [], ...d_n = [])$, where E_i denotes the name of (i)th event, and $d_1 - d_n$ denote the range of concrete values of n data variables mapping to the event E_i . E.g, in our first running example of smart meter, (receive: seg-data=false, command=nil, status=time-out, len (partial)>0) represents the assigned values of selected

¹This is similar to what almost all other invariant detection techniques do, with the exception of DAIKON, which has an automated instrumentation engine.
data variables during *receive* mode. Figure 3.2 shows sample DElog and D|E invariants of motivating example.

We have chosen the above D|E invariant template as a common data invariant template. However, ARTINALI D|E miner is extensible to all templates that Daikon uses for data invariant inference. We avoid using all Daikon templates for three reasons: First, data invariants inferred using various templates are overlapped (e.g, $2 \le X \le 5, Y \ge 5, Y \ge X$). Secondly, the more number of invariants leads to a higher rate of false positives in anomaly detection [6]. Thirdly, CPS has a limited memory capacity which makes using a big IDS model challenging. Thus, a smaller set of rich and stable invariants for CPS IDS model is more desired.



Figure 3.2: Sample DElog and respective D|E invariants with 100% confidence

	E T Invariant Type
Type I	$E_i(t) \rightleftharpoons E_i(t + \frac{1}{Freq_i})$
Type II	$E_j \rightleftharpoons E_i : \Delta t_{ji} \max, \Delta t_{ji} \min$
Type III	$E_i: \Delta t_i \max, \Delta t_i \min$
	D T Invariant Type
Type I	$d_m(T_i \le t \le T_j) = []$
Type II	$d_m = [] \rightleftharpoons d_n = [] : \Delta t_{ji} \max, \Delta t_{ji} \min$

Table 3.1: E|T and D|T Invariant Types

3.3.2 Block 2. ARTINALI E|T Miner

ARTINALI'S E|T Miner infers the E|T invariants in four steps. First, it creates all consecutive event pairs within one trace annotated with their time differences. Second, it groups the pair of events that are labeled with the same pair name. Third, ARTINALI'S E|T Miner looks for the pair-wise events that are observed in the same order within training traces, and calculates their support. Finally, it merges the time variables within each class to calculate the time boundaries of the paired events, as well as the frequency and the average duration of every event execution. The E|T invariants are classified into three types, as shown in Table 3.1. Type I indicates that event E_i is repeated every $\frac{1}{Freq_i}$ seconds. Type II indicates that the pair of events E_i and E_j are repeated in all traces in the same order, and their time difference is bounded within $\Delta t_{ji}max$ and $\Delta t_{ji}min$. Type III indicates the maximum and minimum duration of event E_i . For the example in Section 2.2, invariant send \Rightarrow send :60.2, 59.9sec showing the frequency of send occurrences in the system, and the invariant send \rightleftharpoons receive : 1.2, 1.5 representing the time boundary (between 1.5 and 1.2) as well as the logical ordering of the events (i.e., send before receive), are both examples of E|Tinvariants. We illustrated sample TElog and DurationLOG files that are fed to E|T miner, and the respective E|T mined invariants in Figure 3.3.



Figure 3.3: Sample TElog and DurationLOG files and respective E|T invariants.

3.3.3 Block 3. ARTINALI D|T Miner

According to the formulation described for D|T invariants, ARTINALI combines the outputs of D|E and E|T miners to generate the real-time data invariants (D|T invariants). We define two types of data invariants (Table 3.1), and we explain each type using the example in Section 2.2.

Type I represents the distribution of valid data values of variable d_m within time slot $T_i \leq t \leq T_j$. For example, seg-data $(T_1 \leq t \leq T_2) = true$ means that the only valid value of variable seg-data is *true* during the time interval $T_1 \leq t \leq T_2$. Note that we differ here from Daikon data invariants (e.g. seg-data=true,false), as they only express the valid values of data

invariants without considering the time.

Type II captures the relationship of data invariants between two consecutive events. As explained in previous section, every two consecutive events have a bounded time difference $(T_i + \Delta t_{ji} \min \leq T_j \leq T_i + \Delta t_{ji} \max)$. As a result, the data invariants associated with those events have the same time difference. In other words, data invariant $d_j = []$ holds true *until* data invariant $d_i = []$ becomes true, while $\Delta t_{ji} \max$ and $\Delta t_{ji} \min$ specifies the time difference boundaries between those data invariants. Figure 3.4 shows examples of two types of D|T invariants that ARTINALI D|T Miner generates for our running example. For example, invariant seg-data = $true \rightleftharpoons$ seg-data = false : 1.2, 0.1sec; i.e., seg-data = trueholds true *until* seg-data is assigned value false, in a time interval ranging between 0.1 and 1.2 seconds.



Figure 3.4: Representation of ARTINALI D|T invariants

3.3.4 Block 4. IDS Prototype

As explained in the previous sections, the ARTINALI Miners derive three classes of invariants that comprise the final CPS model. The CPS model not only satisfies the mined invariants but also admits the correct paths within executions. It is used as an input to our IDS prototype for monitoring attacks. Our IDS prototype consists of two components: the *Tracing module* and the *Intrusion detector*. The tracing module is in charge of collecting the required information from the program's execution and logging it. This module is the same as the ARTINALI Logger that instruments the code and collects logs, but with the difference that it is deployed on the production system. The collected information is fed to the intrusion detector, which periodically processes the log file and checks it against the invariants derived from the CPS model.

3.4 Summary

In this chapter, we introduced ARTINALI for mining likely invariants through dynamic analysis in CPS systems, for specification-based IDSes. The main innovation in ARTINALI is that it incorporates time as a first-class notion in the mined invariants, in addition to the traditional data and event invariants. This is important for two reasons. First, most CPSes have real-time constraints, and hence their operational correctness depends on both logical correctness, and correct timing behavior. Hence, incorporating time is essential for detecting many common security attacks in these systems. Secondly, CPSes have predictable timing behaviors to a first order of approximation. and hence leveraging this predictability leads to higher accuracy (i.e., lower false-positives and negatives). However, incorporating time in dynamic invariant detection techniques increases the complexity of the learning due to the much larger state space that needs to be covered. To alleviate this issue, we break up the problem of learning invariants along the three dimensions into problems of learning invariants along two dimensions, namely data-events and events-time, and then combine them into data-events-time invariants. As a proof of concept, we have also built an ARTINALI-based IDS prototype, that is fed by multi-dimensional model generated by ARTI-NALI, and use it in the context of two CPS systems, namely i) advanced metering infrastructures, ii) and smart artificial pancreas (Chapter 4).

Chapter 4

Experimental Setup

This section first presents the details of two CPS platforms as case studies, and then the experimental procedure for evaluating the IDS on the two platforms.

4.1 Case Studies

We chose two CPS platforms as case studies to evaluate the efficacy of the invariants generated by ARTINALI and the other tools. Note that unlike generic applications, there are few publicly available open-source CPS platforms that are also security-critical. Furthermore, there is a significant amount of effort involved in setting up a CPS platform and generating execution traces from it.

4.1.1 Advanced Metering Infrastructure (AMI)

Advanced Metering Infrastructure (AMI) systems are deployed on smart electric power grids. Smart meters are key components of AMI that provide a two-way communication with the utility provider[44]. The large scale deployment of smart meters and the discovery of many vulnerabilities in these systems [49, 57], make them good candidates to evaluate our work. According to Figure 4.1(a), a generic smart meter is composed of two main components, namely the meter and the gateway (controller). The meter component receives power consumption data (PCD) through analog front end sensors, and stores them in the memory. The controller component is the communication bridge between the meter and the utility provider's server, passing server commands to the meter, and sending consumption data back to the server at specific time intervals (more details in [49]).

Our testbed: We use SEGMeter [1], an open source smart meter to evaluate our IDS prototype. SEGMeter is implemented using the Lua language, and consists of 2500 lines of code (excluding libraries).

4.1.2 Smart Artificial Pancreas (SAP)

Diabetic patients are migrating from the traditional glucose meter and manual insulin injection systems to continuous glucose monitoring and autonomous insulin delivery devices [33], which are referred to as Smart Artificial Pancreas (SAP). Since attacks to a SAP can threaten the patient's life, these systems are highly security-critical [41]. Hence, we selected SAP as our second case study to evaluate ARTINALI. The main building blocks of a generic SAP are a Continuous Glucose Monitor (CGM), an insulin pump, and a controller (As illustrated in Figure 4.1(b)). These are commonly connected through a wireless network to form a real-time monitoring and response loop. The CGM samples the patient's blood glucose (BG) levels on a regular basis and sends it to the controller. The insulin pump is a wearable medical device that is used for automatic injection of insulin through subcutaneous infusion. It may deliver insulin in two doses: bolus and *basal*. Each type has specific injection time, rate, and dosage based on the patient's needs. The controller controls the closed loop in the SAP. It receives the measured BG from CGM, and issues the suitable actuation command for correcting the sugar level.

Our testbed: We used Open Artificial Pancreas System (OpenAPS) [32], an open source SAP, as a second use case to evaluate our IDS prototype. OpenAPS implements the controller component of a SAP in JavaScript, and consists of 2000 lines of code (excluding libraries). We simulated a simple CGM and an insulin pump to close the loop, as we did not have access to a patient with a real insulin pump and glucose meter. OpenAPS provides a set of test cases that take different BG values as input and process them for calculating basal rate of insulin, which we use as a baseline for our experiments.



Figure 4.1: High level block diagram of (a) Smart meter and (b) Smart Artificial Pancreas.

4.2 Experimental Procedure

Figure 4.2 shows the overall procedure that we follow. In addition to generating the CPS model using ARTINALI, we generate three other models (invariants sets) using Daikon, Texada, and Perfume for comparison purposes. We downloaded the latest versions of these tools from their respective websites [2–4]. We do not run the instrumentation front-end of Daikon (i.e., Kvasir), as our goal was to generate data invariants based on the event traces we logged. We choose these three tools to represent the first, second and fourth classes of invariants as described in Chapter 2. We do not choose the tools in the third category, namely GK-tail and Quarry, as we use Daikon to find data invariants for the events that we identified in the system. Therefore, the invariants generated by Daikon cover the third class of invariants in our experiments (i.e., D|E invariants).

There are 22 system calls in SEGMeter's code, and 4 system calls in the OpenAPS code. We consider all of them as events. Tables 4.1 and 4.2 present the types of invariants and the number of invariants generated by the three tools and ARTINALI for the SEGMeter and OpenAPS platforms respectively. As can be seen, ARTINALI generates invariants in the Time, D|E, E|T, and D|T categories, while DAIKON, Texada and Perfume only generate invariants in the D|E, Event and E|T categories respectively.

Because the format of the invariants generated by these other tools may be different from that expected by our IDS, we wrote scripts to convert the invariants to be in the format expected by the IDS interface. ARTINALI directly generated invariants in the proper format. In case a tool did not generate a certain kind of invariant (e.g., D|E), we leave that invariant file blank. The generated invariant sets are all fed into the IDS as inputs, and their efficacy is evaluated on different platforms.

We divide the experiment into a training phase and a testing phase for each system. We first obtain execution traces from the two platforms under normal operation, and randomly divide them into a set of training traces (train) and testing traces (test). We then choose different training set sizes for each invariant detection system to optimize the false positive (FP) and false negative (FN) ratios for that system. Finally, we evaluate the FP ratios of the invariants using the *test* traces, and the FN ratios using the attack models described in the next section.



Figure 4.2: Overall experimental process of running the IDS

The IDS is implemented in Python, and consists of about 1000 lines of code. Since the IDS is run on the CPS platform, which is often resource

Table 4.1: Types and number of inferred invariants for SEGMeter across tools

	Event	Time	D E	E T	D T
Daikon	-	-	24	-	-
Texada	158	-	-	-	-
Perfume	-	-	-	158	-
ARTINALI	-	12	24	37	24

Table 4.2: Type and number of inferred invariants for OpenAPS across tools

	Event	Time	D E	E T	D T
Daikon	-	-	22	-	-
Texada	57	-	-	-	-
Perfume	-	-	-	57	-
ARTINALI	-	4	22	18	7

constrained, it is important to minimize its overheads. We measure the IDS's time and space overhead for the SEGMeter platform in Chapter 6-Section 6.3.5 and Section 6.3.4. Because we run the OpenAPS platform in a simulator, as we did not have access to its hardware, we do not measure the IDS overheads on OpenAPS.

Chapter 5

Evaluation Against Targeted Attacks

In this chapter, we discuss the potential targeted attacks and how we derive them for SEGMeter and OpenAPS platforms. We then evaluate the IDS seeded by ARTINALI and other tools against the attacks. Note that we used attack trees based on prior attacks against similar systems to generate the attacks to minimize bias and model realistic attacks. We found that ARTINALI was able to detect all the attacks, while none of the other tools do so. This is because all the attacks involved violations of the interplay among data, events, and time.

5.1 AMI Attacks

Energy fraud is a major class of AMI attacks, and can result in *Power Con*sumption Data (*PCD*) loss and improper billing [36]. We came up with an attack tree for energy fraud in AMI (shown in Figure 5.1), based on attacks introduced in previous work [36, 44, 50, 57]. There are three major branches in this tree, namely i) Measurement tampering, ii) Storage tampering, and iii) Network tampering. Corresponding to each branch, we developed the concrete attack actions as the leaves of the tree as follows.

5.1.1 Synchronization tampering (Blocks A1 - A4)

Synchronization tampering attack occurs due to modification of the time of *send* and *receive* modes in AMI. We found that the communication between the AMI and the server is synchronized by a vulnerable function (*get-data*-





Figure 5.1: Attack tree for AMI

timer()) in the controller unit. The controller frequently checks the time with the sever to decide when to request for data measured by the meter. If a malicious user delays the server commands, the controller will not receive data in the expected time, which leads to data loss, and improper calculation of final PCD.

5.1.2 Meter spoofing (Blocks B1 - B5)

In a smart grid, AMIs communicate with the server using a unique name or ID. The controller unit is able to be connected to more than one meter, collects the PCDs, and send them along with the meter's ID to the server. As the controller cannot differentiate between normal and abnormal messages, it can be tricked by falsified inputs sent by an attacker instead of the meter. This attack is called *meter spoofing attack*. We found that spoofing the meter only requires the meter's ID that is printed on the meter's nameplate.

5.1.3 Message dropping (Blocks C1 - C5)

An attacker may be able to drop the messages (i.e., a part of energy usage) after bypassing the meter and removing the logged PCD history. A simple

Table 5.1: ARTINALI invariants to detect the example attacks in AMI.

Attack	Detecting Invariant			
Synchronization tampering	(1) send $(T0+K \cdot 60) \rightleftharpoons$ send $(T0+(K+1) \cdot 60), \forall k \ge 0$			
Message dropping	(2) recv (T1) \rightleftharpoons recv (T1+1)			
Meter spoofing	(3) node-name(T0+N \cdot 60) = Node B, \forall N \geq 0			

way to mount this attack is to intercept the communication between the meter and the controller, and control what traffic to block and what to pass through (e.g., through a firewall). Hence, the blocked traffic would not be included in PCD calculations.

5.2 Detection of AMI attacks

We ran the ARTINALI-based IDS on the example attacks, and found that it detected all of them. Table 5.1 indicates the important invariants that are derived by ARTINALI, which detect the attacks presented in the previous section.

5.2.1 Synchronization tampering

As synchronization tampering attack modifies the timing of *send* and *receive* operations of SEGMeter, we picked events *send* and *receive* as relevant events to explain this attack. We can see in row 1 of Table 5.1 that the ARTINALI invariant captures the sequence of these events during normal operation, i.e., *send* operation happens every 60 seconds, and *receive* is repeated every 1 second. Thus, this invariant detects the attack as the timing of the events is violated by the attack.

5.2.2 Message dropping

If we assume the attacker drops one or more messages from meter, the dropped messages will not be received at the expected time slots by the controller. As a result, the frequency of receiving messages in controller will change. This attack breaks the invariant number (2) in Table 5.1, which represents the time frequency of receive function which is 1003 milliseconds $(\cong 1sec)$ within one full execution path. Thus this attack is also detected.

5.2.3 Meter spoofing

To detect meter spoofing attack, we selected two receive events (recvA and recvB) from two different meters (node A and node B) that are connected to the same controller, and analyzed the respective invariants. For example, $nodeName(T0+N^*60) = Node B, \forall N \ge 0$ specifies that the valid value of nodeName at T0 + N * 60 is Node B. If the identity of node A is stolen by node B, it sends its messages every 60 seconds under the name nodeA. As a result, variable nodeName attached to event recvB, becomes nodeA. Thus, the invariant number (3) in Table 5.1 is violated.

5.3 SAP Attacks

Diabetic therapy tampering is one of the highest severity threats for patients, as it can result in death or severe health complications. We developed an attack tree for diabetic therapy tampering based on publicly available reports of attacks on SAPs [33, 41], in Figure 5.2. We consider three classes of attacks based on the tree.

5.3.1 CGM spoofing attack (Blocks A1 - A4)

The CGM spoofing attack injects false into the communication channel between CGM and controller making the controller think that the glucose level is either higher or lower than it actually is. There are two ways that CGM spoofing can be accomplished. First, if the sensor data format is unknown, then a replay attack can be used. In this case, a sensor value read in the past can be re-sent (e.g., by a RF module [33]) to the controller. This would cause the controller unit to indicate an outdated glucose level rather than the actual one. Second, if the format of sensor data is known to hacker, she can send the false data at random time intervals to mislead the controller.



Figure 5.2: Attack tree for SAP

Table 5.2: ARTINALI invariants to detect example attacks in OpenAPS.

Attack	Detecting invariant
CGM spoofing	(1) read (t) \rightleftharpoons read (t+5)
Stop basal injection	(2) $(120 \le BG \le 485) \Rightarrow (0.9 \le basal.rate \le 3.5) : 1.99, 0.464$
Resume basal injection	(3) BG $\leq 75 \rightleftharpoons$ basal.rate=0 : 1.99, 0.464

5.3.2 Basal tampering (Blocks B1 - B5)

The basal tampering attack may be accomplished in two different scenarios. The attacker may issue a command for i) stopping the basal injection (e.g., basal.rate = 0) when it is required for patient, or ii) resume the basal injection (basal.rate > 0) when it has to be stopped. These attacks may be mounted using a software radio board that fully controls the SAP [33, 41], and transmits the malicious commands to the pump. To accomplish the attack, the attacker needs to spoof the PIN number of the controller, and the format of transmission packets - both of these can be done by an eavesdropping attack.

5.4 Detection of example attacks in SAP

We mounted the attack examples on the SAP system we considered (i.e., OpenAPS), and found that the ARTINALI-based IDS is able to detect all of them. There are four events in the SAP, namely 1) send(BG) or sending blood glucose by CGM, 2) read(BG) or reading BG by the controller, 3) send(basal.rate) or sending basal rate to pump by the controller, and 4) recv(basal.rate) or receiving basal rate by pump. We used these events as the basis for mining 51 invariants for OpenAPS's IDS model. Due to space constraints, we do not present all inferred invariants, but only those that detect the example attacks (Table 5.2).

5.4.1 CGM spoofing attack

We selected read(BG) in controller as the relevant event, and analyzed the inferred invariants for this event to analyze CGM spoofing attack. Under normal conditions, the transmission of measured Blood Glucose (BG) to CGM occurs at deterministic, periodic times (e.g., every five minutes). This property is represented in our model as time frequency of event read(BG), that is $read(t) \rightleftharpoons read(t+5)$. Using the above property, it would be possible to detect malicious sensor reading from any external source that performs replay attack or transmits wrong data at random time intervals to the controller as the frequency of reading data by controller would change.

5.4.2 Basal tampering attack

As previously explained, the basal tampering attack may be accomplished in two different scenarios: i) stop basal injection (basal.rate = 0) when it is required, and ii) resume basal injection (basal.rate > 0) when it is not required. These attacks break the invariants shown in Table 5.2. The invariant number (2) indicates that if BG is higher than the normal range, the patient needs insulin (i.e., basal.rate > 0). However, the *stop insulin injection* attack makes the *basal.rate* value to be 0, which breaks invariant number (2). Similarly, the invariant number (3) in Table 5.2 shows that for low BG ranges (e.g., BG = 45), the patient does not need insulin (i.e, basal.rate must be 0), but resume basal injection attack sends a command (basal.rate > 0) to the SAP to inject insulin. As a result, invariant number (3) is violated.

5.4.3 Summary

Thus, we see that the ARTINALI-based IDS is able to detect all six attacks that we considered (3 for AMI, and 3 for SAP). On the other hand, none of the other three systems (i.e., DAIKON, Texada, and Perfume) detect even a single one of the attacks. This is because all the attacks involved violations of the interplay among data, events and time

Chapter 6

Evaluation Against Arbitrary Attacks

In the previous chapter, we evaluated the IDS based on ARTINALI on targeted attacks. However, evaluation of security techniques using a small number of targeted (hand-crafted) attacks might not be sufficient for CPS systems for two reasons. First, CPSes are new systems for which there are few real attacks - hence they need protection from zero-day or unknown attacks. This is especially the case for security-critical CPSes such as smart medical devices. Secondly, unlike general computer systems, CPSes can be difficult to upgrade and patch frequently, and hence, they need resilient IDSes against arbitrary attacks. Therefore, in this chapter, we evaluate the efficacy of invariants generated by ARTINALI and the other three invariantdetection techniques for the CPS platforms using arbitrary attacks.

6.1 Arbitrary Attack Model

We use fault injection (i.e., mutation testing) to emulate arbitrary attacks. Fault injection has been used to study the effects of attacks in previous work [49]. Note that these are not complete attacks, but rather form the building blocks of attacks. We deploy different types of mutation in the program's code, as follows.

- *Data mutations*, which change the runtime values of data variables in the code;
- *Branch flipping*, which change the normal execution flow of the program by flipping branch conditions;

• *Artificial delay insertion*, which modify the normal timing behavior of the program.

Each of the above categories emulate different security issues. By performing data mutations, an attacker can change critical data in the program to their advantage. Such attacks can be accomplished by exploiting memory corruption vulnerabilities or race conditions in the program. For instance, [48] investigated an attack to smart facial recognition systems caused by exploiting a mis-classification bug (CVE-2016-1516) in the controller algorithm through input data mutation. Likewise, branch flipping can lead to illegitimate control flow paths being taken in the program, to accomplish the attacker's ends. Such attacks can occur due to code injection or semantic vulnerabilities. As an example, [10] indicated that how attacker is able to change the sequence of actions in a smart car through exploiting the buffer overflow vulnerability in the telematics of car. Finally, artificial delays can allow attackers to change the timing of the system's actions, and delay essential functions, or cause other functionality to be suppressed, again to their advantage. Synchronization tampering attack is one example of such attacks [36]. Through these mutations, we can emulate a wide variety of attacks, without a predefined target, thus avoiding bias and allowing modeling of hitherto unknown attacks.

We performed 156 and 125 code mutations for SEGMeter and OpenAPS respectively. We manually seeded each of these mutations in the source code of the respective systems, by randomly sampling the corresponding program points in the program's code. While this could have been automated by a fault injection tool (e.g., LLFI [5]), the languages in which the two systems were implemented, JavaScript and Lua, were not supported by existing tools. So we had to perform mutations manually. However, we attempted to choose the program points randomly before performing the experiment to avoid biasing our evaluation.

After mutating the code, we can observe one of four outcomes.

- *Crash*, in which the program is aborted (exception);
- *Hang*, in which the program goes into an infinite loop or deadlocks;



Figure 6.1: Fault injection impact(%): (a) data mutation in SEGMeter, (b) branch flip in SEGMeter, (c) artificial delay in SEGMeter, (d) data mutation in OpenAPS, (e) branch flip in OpenAPS, (f) artificial delay in OpenAPS

- *SDC (Silent Data Corruption)*, in which the outcome of the program is different from a fault-free execution;
- *No corruption*, in which the outcome of the program does not show any observable impact with respect to fault masking or non-triggering faults. Internal states might however be corrupted.

Note that in the context of this study, we are interested only in *SDC* and *No corruption* outcomes, as the *Crash* and *Hang* outcomes can easily be detected without an IDS. Therefore, we need an IDS for the *SDC* and *No corruption* outcomes which comprise about 85% of the outcomes on average (according to Figure 6.1).

6.2 Evaluation Metrics

Accuracy: We use three metrics to measure the effectiveness of our IDS from the accuracy point of view.

- False Negative ratio (FN), which is the ratio of attacks that were undetected by the IDS to the total number of attacks;
- False Positive ratio (FP), which is the ratio of execution traces that were (incorrectly) reported as attacks to the total number of normal traces;
- F-Score(β), which is a computation of the harmonic mean of the true positive ratio (TP), FP and FN².

The variations of the argument β in F-Score(β) allow us to weigh the above metrics differently [46], and obtain different trade-off between FP and FN ratios based on the system requirements. A value of $\beta > 1$ weighs FNs higher, while a value of $\beta < 1$ weighs FPs higher. A value of $\beta = 1$ weighs them both equally. We hypothesize that FPs are more important in smart meters, as a false-alarm leads to added cost to the utility provider who needs to deploy service personnel to investigate the false alarm. An occasional FN may be acceptable in smart meters as the consequence is only a loss of revenue. In the OpenAPS, on the other hand, even a single FN can be fatal to the patient, while a FP may be acceptable if there are other checks in place to filter out FPs (e.g., patient intervention). Hence, for SEGMeter, we select F-Score(0.5), and for OpenAPS, we choose F-score(2) as our reference metric.

Overheads: In addition to the accuracy, we also measure the memory and performance overheads of the IDS.

Memory overhead is defined as the actual memory usage of the IDS. It depends on the size of IDS, the number of invariants that account for the CPS model, and the complexity of invariants (e.g., the invariant $E_j \rightleftharpoons E_i$: Δt_{ji} max, Δt_{ji} min carries more information than the invariant $E_j \rightleftharpoons E_i$, and is hence more complex).

Performance overhead is the increase in execution time as a result of running the CPS on the target platform. This metric reflects the overhead of both the *tracing module* and the *intrusion detector*. Since CPSes run

 $^{{}^{2}}F - Score(\beta) = \frac{(1+\beta^{2})\times TP}{(1+\beta^{2})\times TP + \beta^{2}\times FN + FP}$

continuously for long periods of time, we measure the performance overhead per cycle, where a cycle refers to one full execution of the main loop of the CPS (both the SEGmeter and OpenAPS consist of a single main loop that runs continuously).

6.3 Research Questions (RQs)

We ask the following RQs corresponding to the evaluation metrics.

- **RQ1.** How do we choose the training set size to obtain the best F-Score(β) for each tool?
- **RQ2.** What is the FN ratio incurred by the IDS using the invariants derived by ARTINALI and the other tools ?
- **RQ3.** What is the FP ratio incurred by the IDS using the invariants derived by ARTINALI and the other tools?
- **RQ4.** What is the memory overhead of the IDS when using the invariants derived by ARTINALI and the other tools?
- **RQ5.** What is the performance overhead of the IDS when using the invariants derived by ARTINALI and the other tools?

6.3.1 RQ1. F-Score

As mentioned in Chapter 4, we obtain two sets of traces from each system, namely *train* and *test*. In this RQ, we ask what should be the optimal training set size for each system in order to maximize the corresponding F-Score values. To answer this question, we obtain a total of 40 training traces, and 50 test traces for each system. We then vary the training set size from 5 to 40, in increments of 5. We then run each of the invariant detection tools including ARTINALI on the same training set to derive invariants. We then measure the FP, FN, and F-Score values (0.5, 1, 2) for each invariant detection tool and system, as a function of the training set size.

Figures 6.2, 6.3, 6.4 and 6.5 show the distribution of the amount of false positives (FP), false negatives (FN) and the F-Score computed with $\beta = 0.5, 1, 2$ in relation to the amount of training traces for *SEGMeter*, corresponding to each of the four invariant detection tools, including ARTINALI. Similarly, Figures 6.6, 6.7, 6.8 and 6.9 show the distribution of the amount of false positives (FP), false negatives (FN) and the F-Score computed with $\beta = 0.5, 1, 2$ in relation to the amount of training traces for *OpenAPS*, corresponding to each of the four invariant detection tools. As expected, as the amount of training traces increases, the FP ratio decreases, since a broader set of invariants are extracted; thus a lower amount of legitimate actions are flagged as potential attacks. A consequence is that more attacks are undetected (FN increases), as a more restricted set of invariants can lead to some attacks being undetected. Overall, an increase in the amount of training traces lead to an increase of the F-Score at first, then it stabilizes, at which point an *optimal* amount of training traces have been found (for a given values of β).

Tables 6.1 and 6.2 show the optimal amount of training traces (optimal F-Score) for each invariant detection tool, for SEGMeter and OpenAPS respectively. Recall that we choose F-Score(0.5) for SEGMeter and F-Score(2) for OpenAPS, and hence these are the F-score values we choose for the optimal number of traces. For example, in SEGMeter, a training set size of 20 results in the maximum value of the F-Score(0.5) value of ARTINALI, whereas for OpenAPS, a training set size of 15 results in the maximum value of F-Score(0.5). Likewise, we compute the optimal training set sizes for the three other tools on both platforms. These are the values of the training set sizes we use for deriving the invariants for each tool in the rest of this section. In other words, we find the best configuration of each tool on each platform, and generate invariants using this configuration for comparing the corresponding IDSes.

Table 6.1: Optimal training set size for maximum F-Score(0.5) for SEGMeter across tools, and the corresponding FP and FN ratios.

	Daikon	Texada	Perfume	ARTINALI
F-Score(0.5)	0.721	0.78	0.813	0.898
Num of traces	30	30	35	20
FP (%)	23	15	15	12
FN (%)	57	60	38	2.3

Table 6.2: Optimal training set size for maximum F-Score(2) for OpenAPS across tools, and the corresponding FP and FN ratios.

	Daikon	Texada	Perfume	ARTINALI
F-Score(2)	0.604	0.62	0.686	0.952
Num of traces	30	20	15	15
FP (%)	21	16	22	13.5
FN (%)	61	61	39	2



Figure 6.2: FN, FP and F-Score variations based on number of training traces for SEGMeter's IDS seeded by Daikon invariants.



Figure 6.3: FN, FP and F-Score variations based on number of training traces for SEGMeter's IDS seeded by Texada invariants.



Figure 6.4: FN, FP and F-Score variations based on number of training traces for SEGMeter's IDS seeded by Perfume invariants.



Figure 6.5: FN, FP and F-Score variations based on number of training traces for SEGMeter's IDS seeded by ARTINALI invariants.



Figure 6.6: FN, FP and F-Score variations based on number of training traces for OpenAPS's IDS seeded by Daikon invariants.



Figure 6.7: FN, FP and F-Score variations based on number of training traces for OpenAPS's IDS seeded by Texada invariants.



Figure 6.8: FN, FP and F-Score variations based on number of training traces for OpenAPS's IDS seeded by Perfume invariants.



Figure 6.9: FN, FP and F-Score variations based on number of training traces for OpenAPS's IDS seeded by ARTINALI invariants.

6.3.2 RQ2. False Negatives

In this section, we compare the variation in the FN ratio incurred by the IDS, using invariants extracted by ARTINALI and the other tools. Tables 6.1 and 6.2 also show the FN ratios for each tool for the SEGMeter and OpenAPS systems respectively. We observe that overall, ARTINALI was able to detect around 97.5% of attacks, which means it has an average FN ratio of 2.5%. In contrast, in Perfume, Texada and Daikon the FN ratio was respectively 38.5%, 60.5% and 59% on average, across the two platforms. Thus, the ARTINALI-based IDS reduces the ratio of false negatives by 89 to 95% over other dynamic invariant detection tools.

Figure 6.10 and Figure 6.11 illustrate the FN ratio of the IDS for the three category of attacks (code mutations), as well as the aggregated FN ratio, for each tool, in both SEGMeter and OpenAPS. We discuss the FN ratio for each attack category below:

Data mutations: ARTINALI exhibits the lowest FN rate for data mutations (2 to 3%). This is followed by Daikon, which provides a much lower FN ratio in *data mutation* attacks (15% in SEGMEeter and 17% in Ope-



Figure 6.10: FN(%) of IDS for SEGMeter for different attack types across the tools. Error bars are shown for the 95% confidence interval.



Figure 6.11: FN(%) of IDS for OpenAPS for different attack types across the tools. Error bars are shown for the 95% confidence interval.

nAPS) than Perfume (53% in SEGMEeter and 78% in OpenAPS) and Texada (52% in SEGMEeter and 87% in OpenAPS). This is because DAIKON focuses on data invariants, while Texada and Perfume do not include data invariants in their model. However, the Daikon data model does not include other properties like ARTINALI does, resulting in much higher FNs than ARTINALI.

Branch flipping: Among the other three tools, ARTINALI has the lowest FN rate for branch flipping attacks (1%). Perfume, Texada and ARTINALI exhibit a lower FN ratio compared to Daikon for branch flipping attacks. As these attacks impact the order and sequence of the events in an execution instance, and Daikon does not have event invariants, it shows less sensitivity.

Artificial delay: Again, ARTINALI has a much lower FN ratio (2-3%) than all three tools for artificial delay attacks, followed by Perfume. This is because they both include time in their model. Nevertheless, Daikon and Texada are still able to detect attacks that impact data variables or alter the execution flow of the program.

Overall, the results support our hypothesis that a more comprehensive invariant model, such as ARTINALI, which can find invariants and their constraints along three dimensions, can detect a significantly larger amount of attacks (and hence has fewer FNs).

6.3.3 RQ3. False Positives

In this section, we compare the FP ratio incurred by our IDS when using the invariants derived by ARTINALI against the invariants generated by the other tools (Daikon, Perfume and Texada). The results are shown in Table 6.1 and 6.2 for the SEGMeter and OpenAPS systems respectively. We can observe that in both CPSes, the use of the ARTINALI-generated invariants lead to significantly less false positives compared to the invariants generated by the other tools. More precisely, ARTINALI provides a 20% to 48% improvement of the FP ratio for SEGMeter, and a 16% to 39% improvement of the FP ratio for OpenCPS, over the other tools. These results can be explained by the fact that ARTINALI leverages the correlations among data, event and time dimensions during correct system behavior to generate more stable invariants. ARTINALI infers event invariants that precisely describe the ordering of events in a sequence within an execution flow, and then associates data and time constraints to the events within every path (D|E and E|T). Therefore, during normal operation, the system is unlikely to follow the same path with different associated data and time values in a given execution, which in turn, reduces the probability of false positives. Although the IDS uses the same traces for all tools, none of these tools other than ARTINALI look at the relational constraints of both data and time along the events' paths, resulting in a higher ratio of false positives.

While the FP ratio for the ARTINALI-based IDS is lower than the other tools, it is still high for both platforms. To reduce the FP ratio, one can deploy multiple variants of the code and switch to a different variant when an attack is detected. If the invariant is not violated in the second version, it may be a false positive. Another solution is to remove invariants that exhibit high FP ratios [6], but this may also increase the FN ratio. We defer a detailed treatment of FP ratio reduction to future work.

6.3.4 RQ4. Memory Overhead

We measured the memory consumption of our IDS running on the SEG-Meter platform, using the invariants generated by different tools. We also calculated the number of invariants that ARTINALI and the other tools inferred for both platforms. Our results are shown in Table 6.3 ("Memory usage" row). Generally, invariants that involve two or more dimensions (e.g., E|T invariants) carry more information than the invariants of one dimension (e.g., event invariants), and hence are more complex. We observe that the memory usage grows as the number and complexity of invariants increases. For example, the IDS consumes the maximum memory usage (3.94 MB) when it uses the Perfume-generated invariants, which straddle two dimensions, and have the maximum number of invariants (158 - Table 4.1).

Memory usage (MB)	Daikon 1.24	Texada 3.21	Perfume 3.94	ARTINALI 2.96
Tracing overhead(%) Detector overhead(%) Overall overhead(%)	$22.6 \\ 4.7 \\ 27.3$	$ 13.4 \\ 10.3 \\ 23.7 $	$ 18.8 \\ 13.3 \\ 32.08 $	$23.3 \\ 8.3 \\ 31.6$
Full cycle execution(s) IDS execution time(s)	$60.94 \\ 16.63$		$60.94 \\ 19.57$	

Table 6.3: Memory and performance overhead of IDS, seeded by ARTINALI and the other tools, running on SEGMeter.

Overall, we find that the memory consumption of the IDS with ARTINALIgenerated invariants is lower than those with Perfume or Texada-generated invariants, but higher than those with Daikon-generated invariants. However, the memory usage for all tools is much lower than the available memory in SEGMeter (16 MB).

6.3.5 RQ5. Performance Overhead

In this section, we discuss the performance overhead of our IDS running on the SEGMeter platform, which consists of an embedded microcontroller (Broadcom BCM3302 V2.9 240MHz CPU and 16 MB RAM) running Linux.

Recall that the IDS consists of two components, namely tracing module and intrusion detector module. Table 6.3 (middle part) shows the overheads of the two modules separately for each tool. Each of these measurements is an average of the overhead of 10 execution traces for each tool, where an execution trace is defined as one complete execution of the meter's main loop. We find that ARTINALI and Perfume have the highest aggregate overhead, followed by Daikon, and then Texada. The difference in the overhead is due to the difference in the tracing module, which needs to collect both event and data/time information for ARTINALI and Perfume, compared with Texada (events only), and Daikon (data only).

In addition to the performance overheads, the IDS execution time should be lower than the execution time of the system's cycle, or else it will be

6.4. Summary

unable to keep up with the system. We measure the raw execution times of a full cycle in Table 6.3 (last part). As can be seen from the table, the entire cycle takes about 60 seconds (1 minute). However, the execution of the IDS for each tool takes less than 20 seconds even in the worst case (for Perfume), which is only a third of execution time of the full cycle. Therefore, the IDS is not a bottleneck in any of the four systems, and is easily able to keep up with the system.

Note that the invariant mining process takes place offline, and hence does not contribute to the performance overhead of the IDS running on the CPS platform. Nonetheless, we measured the time to mine invariants using ARTINALI, on a standard desktop system (Intel core i7 processor with 32 GB RAM, running Linux). We found that the time ranges from 8 to 96 seconds in SEGMeter, and from 6 to 36 seconds in OpenAPS. This is a very reasonable cost in most systems. Though this overhead may be higher for larger systems, invariant mining is a one-time process and needs to be redone only when the code is updated.

6.4 Summary

In this chapter, we evaluated the IDS prototype seeded by ARTINALI and other three invariant-detection techniques for two CPS platforms using arbitrary attacks and using accuracy and overhead metrics. We used a modelbased fault injection technique to emulate the building blocks of real attacks. Overall, the results support our hypothesis that a more comprehensive invariant model, such as ARTINALI, which can find invariants and their constraints along three dimensions, can detect a significantly larger amount of attacks (and hence has fewer FNs). Moreover, we observed that in both CPSes, the use of the ARTINALI-generated invariants lead to significantly less false positives compared to the invariants generated by the other tools. We have also found that the memory consumption and performance overhead of the IDS with ARTINALI-generated invariants is lower than those with Perfume or Texada-generated invariants, but higher than those with Daikon-generated invariants. However, the memory usage for all tools is much lower than the available memory in SEGMeter. We also measured the raw execution times of a full cycle of the IDS for each tool, and observed that it takes less than 20 seconds even in the worst case (for Perfume), which is only a third of execution time of the full cycle (60 seconds). Therefore, the IDS is not a bottleneck in any of the four systems, and is easily able to keep up with the system.

Chapter 7

Discussion

In this chapter, we first examine the threats to the validity of our experiments, followed by reflections on ARTINALIS generalizability.

7.1 Threats to Validity

An external threat to the validity is the limited number of CPS platforms considered (two). However, as we have mentioned, finding CPS platforms that are publicly available and security critical is a challenge. We have attempted to mitigate this threat by choosing two fairly diverse platforms, with various time constraints, and different IDS optimization goals. We acknowledge that these platforms exhibit somewhat simple behaviors - however, many CPSes fall into this category [15].

An internal threat to validity is in our evaluation of the efficacy of the invariants for attack detection through fault injection experiments. While not necessarily representative of all security attacks, fault injection allows us to emulate the behavior of potential attackers without biasing the evaluation towards known vulnerabilities (at the time of the evaluation). We have attempted to mitigate this threat by using mutation operators that were used for emulating attacks in prior work [5].

Another external threat to validity is that the IDSes based on specification mining techniques are vulnerable to malicious traces during training/test phases, and ARTINALI is not an exception. We have attempted to reduce the threat in the training phase by inferring the invariants offline. During the testing phase, our IDS and CPS are executed on separate processes, and hence the intrusion detector module of IDS is isolated from the attacks that may compromise the CPS itself. However, protecting the tracing module of IDS (that is located on CPS) against attacks is a challenging problem. We assumed that the IDS is the root of trust for now, but in the future, we plan to deploy trusted computing base strategies for the IDS (e.g., secure enclaves), so that the attacker cannot attack the IDS directly.

Finally, a construct threat to validity is the evaluation metrics used for measuring efficacy. FP and FN ratios have however been used in a lot of prior work on intrusion detection, as have F-scores, and hence we do not believe this is a significant threat. Another potential construct threat is the choice of tools we use for comparing with ARTINALI, but we mitigated this to an extent by first systematically classifying the space of invariant detection techniques, and then choosing the tools in each category.

7.2 Generalizability of ARTINALI

ARTINALI relies upon two features, namely event-based semantics, and conditional independence of time and data (Section). Events are operations that involve interaction with the outside world. Event-based semantics implies that every event takes place in a unique time frame, and hence, there is no concurrency among event executions. Secondly, ARTINALI assumes an event occurs at a specific time interval, and subsequently, data variables are assigned to specific values. Thus, the time and data corresponding to a particular event, are conditionally independent regardless of the dependency among events. These two features are a common paradigm for CPSes, and hence ARTINALI can be generalized to other CPS platforms such as pacemakers and unmanned aerial vehicles.

However, ARTINALI is not applicable to non-CPS platforms for two reasons. First, the non-concurrency of events does not hold in non-CPS platforms such as mobile phones. Secondly, CPS events have limited functionality, and hence inferring invariants for each event is straightforward. Unlike CPSes, in general realtime systems, tasks can be of unbounded complexity. Furthermore, general purpose computers with full preemptive (i.e., non-realtime) operating systems have a large space of potential behaviors, which makes it challenging to learn invariants for them.
Chapter 8

Conclusions and Future Work

8.1 Summary

Cyber-physical systems (CPSes) are becoming increasingly subject to security attacks due to their interconnectedness and relative lack of protection. In this thesis, we attempt to use dynamic invariant detection techniques to build intrusion detection systems for CPSes. Our key insight is that time is a first class constraint in CPS systems, and hence we incorporate time into the invariants, in addition to data and events. We devise an efficient algorithm for learning invariants over the three dimensions of data, events and time, and implement it in a tool called ARTINALI. We demonstrate the use of ARTINALI on two CPS platforms for intrusion detection. We find that ARTINALI has significantly lower false negatives and false positives than other dynamic invariant detection tools, while incurring comparable performance and memory overheads.

The most important aspect of this work is providing insights regarding overcoming CPS constraints such as real-time constraints and resource constraints, for security developers. Many security solutions rely on a model for the correct behavior of the system they are monitoring. We hypothesized that a more comprehensive model leads to a higher coverage against security attacks. Our results support this hypothesis for ARTINALI, which incorporates real-time constraints along three dimensions into the model, and hence is able to detect a significantly larger amount of attacks.

On the other hand, the complete system model may be large, and check-

8.2. Future work

ing it requires more resources. Unlike the other specification mining techniques (such as Daikon that infers the relationship among all data variables at the entry and exit point of all method calls in the program), ARTI-NALI only captures system calls (which are located on the attack surface) as events to generate the CPS model for attack detection, and hence does not model those parts of the code that are not security-critical. Our result shows that even a selective model, that requires significantly less resources than the complete model, may provide reasonable security coverage. In other words, approximate security solutions are effective for CPSes and provide high enough coverage while requiring significantly less resources. Such techniques may be integrated with CPS software and hence, make it easier to build security solutions that meet the CPS requirements.

In addition, ARTINALI works at *event granularity* to mine invariants. In spite the fact that ARTINALI captures all system calls as events for security purposes, in ARTINALI, events are *user-defined*. This flexibility enables users to optionally customize the level of granularity by choosing another type of events, or prune the space of events by specifying only the *important* ones based on the system's requirements.

Moreover, ARTINALI is a specification mining technique that has been primarily developed for CPS security space. However, as it generates a multi-dimensional model including many rich properties of correct real-time behavior of a program, it is applicable in a broader context other than security. More particularly, it can be used for software *reliability* purposes, in a wide assortment of tasks, such as non-malicious bug detection, software testing, data structure repair, and debugging of applications. Furthermore, ARTINALI-generated invariants can be useful for supporting program evolution and comprehension as well.

8.2 Future work

There are three potential directions in which this thesis can be extended in the future.

8.2.1 Extending the Generalizability of ARTINALI

In this thesis, we focused on two families of CPS platforms including smart meters and smart artificial pancreases as case studies. Although we expect the same technique to apply to similar CPS platforms, it would be interesting to examine the generalizability of ARTINALI to more complex CPS platforms such as unmanned aerial vehicles (i.e., drones). We define complexity as the total number of source lines of code and having more diverse functionalities (i.e., multiple operational modes). Prior work [23, 47] has shown there are significant challenges in modeling the behavior of drones. For instance, drones operate in various flight modes with different level of autonomy, and in non-autonomous modes, they show more uncertainty not only in the timing behavior but also in the functional behavior. Hence, generating system model for non-autonomous modes in these systems is a challenge. Moreover, drones behave differently in various flight states (including landing, taking off, hovering, etc). Therefore, designing an CPS model for the entire system that reflects the properties of each state precisely is challenging.

8.2.2 Optimizing ARTINALI for the Scalability

In this thesis, we found that ARTINALI-based IDSes impose acceptable overheads on our CPS platforms. However, the resource constraints of a CPS platform with larger code size may degrade the scalability of our technique as the resource utilization of IDS increases. According to our experimental results, a large fraction of IDS overhead is due to the tracing module collecting online information about the system. Prior work [52] has found there is a tradeoff between attack coverage and the amount of data collected during tracing. Hence, a potential future direction is to optimally deploy the tracing with respect to IDS goals including *coverage* and *scalability*, and CPS constraints including *memory usage*, *performance overhead* and *computational power overhead*.

8.2.3 Extending ARTINALI for Attack Diagnosis

In this thesis, we proposed a multi-dimensional invariant mining technique for attack detection in CPSes. However, once an attack is detected, it needs to be mitigated, and attack diagnosis bridges the gap between attack detection and mitigation. Attack diagnosis is a promising approach to carrying out a propagation analysis from source of attacks (cyber vulnerabilities) to the corresponding effects on the physical system [56]. Such analysis is useful to rank the scope and severity of the cyber attacks and thereby, prioritize among various mitigation mechanisms to be selected with respect to CPS constraints. A potential future direction is to incorporate the dynamic invariants generated by ARTINALI for attack diagnosis and mitigation.

- Smart energy groups home page. http://smartenergygroups.com., 2011.
- [2] Perfume User Manual. http://people.cs.umass.edu/~ohmann/ perfume/, 2014.
- [3] Texada User Manual. https://bitbucket.org/bestchai/texada/, 2016.
- [4] The Daikon Invariant Detector User Manual. https://plse.cs. washington.edu/daikon/download/doc/daikon.html, 2017.
- [5] Maryam Raiyat Aliabadi and Karthik Pattabiraman. Fidl: A fault injection description language for compiler-based sfi tools. In *International Conference on Computer Safety, Reliability, and Security*, pages 12–23. Springer, 2016.
- [6] Leonardo Aniello, Claudio Ciccotelli, Marcello Cinque, Flavio Frattini, Leonardo Querzoni, and Stefano Russo. Automatic invariant selection for online anomaly detection. In *International Conference on Computer* Safety, Reliability, and Security, pages 172–183. Springer, 2016.
- [7] Robin Berthier, William H Sanders, and Himanshu Khurana. Intrusion detection for advanced metering infrastructures: Requirements and architectural directions. In Smart Grid Communications (SmartGrid-Comm), 2010 First IEEE International Conference on, pages 350–355. IEEE, 2010.
- [8] Alvaro Cardenas, Saurabh Amin, Bruno Sinopoli, Annarita Giani, Adrian Perrig, and Shankar Sastry. Challenges for securing cyber physi-

cal systems. In Workshop on future directions in cyber-physical systems security, page 5, 2009.

- [9] Alvaro A Cardenas, Saurabh Amin, and Shankar Sastry. Secure control: Towards survivable cyber-physical systems. In *Distributed Computing* Systems Workshops, 2008. ICDCS'08. 28th International Conference on, pages 495–500. IEEE, 2008.
- [10] Stephen Checkoway, Damon McCoy, Brian Kantor, Danny Anderson, Hovav Shacham, Stefan Savage, Karl Koscher, Alexei Czeskis, Franziska Roesner, Tadayoshi Kohno, et al. Comprehensive experimental analyses of automotive attack surfaces. In USENIX Security Symposium. San Francisco, 2011.
- [11] E Chien, N Falliere, and LO Murchu. W32. stuxnet dossier. symantec security response, 2010.
- [12] Mihai Christodorescu, Somesh Jha, and Christopher Kruegel. Mining specifications of malicious behavior. In Proceedings of the the 6th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on The foundations of software engineering, pages 5–14. ACM, 2007.
- [13] Christoph Csallner, Nikolai Tillmann, and Yannis Smaragdakis. Dysy: Dynamic symbolic execution for invariant inference. In *Proceedings of the 30th international conference on Software engineering*, pages 281–290. ACM, 2008.
- [14] Barthélémy Dagenais and Martin P Robillard. Creating and evolving developer documentation: understanding the decisions of open source contributors. In Proceedings of the eighteenth ACM SIGSOFT international symposium on Foundations of software engineering, pages 127– 136. ACM, 2010.
- [15] Patricia Derler, Edward A Lee, and Alberto Sangiovanni Vincentelli. Modeling cyber-physical systems. *Proceedings of the IEEE*, 100(1):13–28, 2012.

- [16] John Eidson, Edward A Lee, Slobodan Matic, Sanjit A Seshia, and Jia Zou. A time-centric model for cyber-physical applications. In Workshop on Model Based Architecting and Construction of Embedded Systems (ACES-MB), pages 21–35, 2010.
- [17] Michael D Ernst, Jake Cockrell, William G Griswold, and David Notkin. Dynamically discovering likely program invariants to support program evolution. *IEEE Transactions on Software Engineering*, 27(2):99–123, 2001.
- [18] Earlence Fernandes, Jaeyeon Jung, and Atul Prakash. Security analysis of emerging smart home applications. In Security and Privacy (SP), 2016 IEEE Symposium on, pages 636–654. IEEE, 2016.
- [19] Mark Gabel and Zhendong Su. Symbolic mining of temporal specifications. In Proceedings of the 30th international conference on Software engineering, pages 51–60. ACM, 2008.
- [20] Pedro Garcia-Teodoro, J Diaz-Verdejo, Gabriel Maciá-Fernández, and Enrique Vázquez. Anomaly-based network intrusion detection: Techniques, systems and challenges. *computers & security*, 28(1):18–28, 2009.
- [21] Dan Goodin. US spy drone hijacked with GPS spoof hack, report says. http://www.theregister.co.uk/2011/12/15/us_spy_drone_ gps_spoofing/, 2011. [Online; accessed 15-Dec-2011].
- [22] Gösta Grahne and Jianfei Zhu. Fast algorithms for frequent itemset mining using fp-trees. *IEEE transactions on knowledge and data engineering*, 17(10):1347–1362, 2005.
- [23] Hengle Jiang, Sebastian Elbaum, and Carrick Detweiler. Reducing failure rates of robotic systems though inferred invariants monitoring. In Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on, pages 1899–1906. IEEE, 2013.

- [24] Sami Kamkar. SkyJack. https://samy.pl/skyjack/, 2013. [Online; accessed 19-Dec-2013].
- [25] Hermann Kopetz and Günther Bauer. The time-triggered architecture. Proceedings of the IEEE, 91(1):112–126, 2003.
- [26] Karl Koscher, Alexei Czeskis, Franziska Roesner, Shwetak Patel, Tadayoshi Kohno, Stephen Checkoway, Damon McCoy, Brian Kantor, Danny Anderson, Hovav Shacham, et al. Experimental security analysis of a modern automobile. In 2010 IEEE Symposium on Security and Privacy, pages 447–462. IEEE, 2010.
- [27] Ted Kremenek, Paul Twohey, Godmar Back, Andrew Ng, and Dawson Engler. From uncertainty to belief: Inferring the specification within. In Proceedings of the 7th symposium on Operating systems design and implementation, pages 161–176. USENIX Association, 2006.
- [28] Cheolhyeon Kwon, Weiyi Liu, and Inseok Hwang. Security analysis for cyber-physical systems against stealthy deception attacks. In American Control Conference (ACC), 2013, pages 3344–3349. IEEE, 2013.
- [29] Neal Leavitt. Researchers fight to keep implanted medical devices safe from hackers. *Computer*, 43(8):11–14, 2010.
- [30] Caroline Lemieux. Mining temporal properties of data invariants. In 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, volume 2, pages 751–753. IEEE, 2015.
- [31] Caroline Lemieux, Dennis Park, and Ivan Beschastnikh. General ltl specification mining (t). In Automated Software Engineering (ASE), 2015 30th IEEE/ACM International Conference on, pages 81–92. IEEE, 2015.
- [32] Dana Lewis. Introducing the # openaps project. 2015.
- [33] Chunxiao Li, Anand Raghunathan, and Niraj K Jha. Hijacking an insulin pump: Security attacks and defenses for a diabetes therapy

system. In e-Health Networking Applications and Services (Healthcom), 2011 13th IEEE International Conference on, pages 150–156. IEEE, 2011.

- [34] Hui Lin, Homa Alemzadeh, Daniel Chen, Zbigniew Kalbarczyk, and Ravishankar K Iyer. Safety-critical cyber-physical attacks: Analysis, detection, and mitigation. In *Proceedings of the Symposium and Boot*camp on the Science of Security, pages 82–89. ACM, 2016.
- [35] Davide Lorenzoli, Leonardo Mariani, and Mauro Pezzè. Automatic generation of software behavioral models. In *Proceedings of the 30th international conference on Software engineering*, pages 501–510. ACM, 2008.
- [36] Stephen McLaughlin, Dmitry Podkuiko, Sergei Miadzvezhanka, Adam Delozier, and Patrick McDaniel. Multi-vendor penetration testing in the advanced metering infrastructure. In *Proceedings of the 26th Annual Computer Security Applications Conference*, pages 107–116. ACM, 2010.
- [37] Robert Mitchell and Ing-Ray Chen. A survey of intrusion detection techniques for cyber-physical systems. ACM Computing Surveys (CSUR), 46(4):55, 2014.
- [38] Robert Mitchell and Ing-Ray Chen. A survey of intrusion detection techniques for cyber-physical systems. ACM Computing Surveys (CSUR), 46(4):55, 2014.
- [39] Gail C Murphy, David Notkin, and Kevin Sullivan. Software reflexion models: Bridging the gap between source and high-level models. ACM SIGSOFT Software Engineering Notes, 20(4):18–28, 1995.
- [40] Tony Ohmann, Michael Herzberg, Sebastian Fiss, Armand Halbert, Marc Palyart, Ivan Beschastnikh, and Yuriy Brun. Behavioral resourceaware model inference. In Proceedings of the 29th ACM/IEEE international conference on Automated software engineering, pages 19–30. ACM, 2014.

- [41] Jerome Radcliffe. Hacking medical devices for fun and insulin: Breaking the human scada system. In *Black Hat Conference presentation slides*, volume 2011, 2011.
- [42] Orna Raz, Philip Koopman, and Mary Shaw. Semantic anomaly detection in online data sources. In Software Engineering, 2002. ICSE 2002. Proceedings of the 24rd International Conference on, pages 302–312. IEEE, 2002.
- [43] Fred Samland, Jana Fruth, Mario Hildebrandt, Tobias Hoppe, and Jana Dittmann. Ar. drone: security threat analysis and exemplary attack to track persons. In *Proceedings of the SPIE*, volume 8301, 2012.
- [44] Florian Skopik, Zhendong Ma, Thomas Bleier, and Helmut Grüneis. A survey on threats and vulnerabilities in smart metering infrastructures. *International Journal of Smart Grid and Clean Energy*, 1(1):22– 28, 2012.
- [45] Sean W Smith. Security and privacy challenges in the smart grid. 2009.
- [46] Marina Sokolova, Nathalie Japkowicz, and Stan Szpakowicz. Beyond accuracy, f-score and roc: a family of discriminant measures for performance evaluation. In Australasian Joint Conference on Artificial Intelligence, pages 1015–1021. Springer, 2006.
- [47] Yunmok Son, Hocheol Shin, Dongkwan Kim, Young-Seok Park, Juhwan Noh, Kibum Choi, Jungwoo Choi, Yongdae Kim, et al. Rocking drones with intentional sound noise on gyroscopic sensors. In USENIX Security, pages 881–896, 2015.
- [48] Rock Stevens, Octavian Suciu, Andrew Ruef, Sanghyun Hong, Michael Hicks, and Tudor Dumitraş. Summoning demons: The pursuit of exploitable bugs in machine learning. arXiv preprint arXiv:1701.04739, 2017.
- [49] Farid Molazem Tabrizi and Karthik Pattabiraman. Flexible intrusion detection systems for memory-constrained embedded systems. In

Dependable Computing Conference (EDCC), 2015 Eleventh European, pages 1–12. IEEE, 2015.

- [50] Farid Molazem Tabrizi and Karthik Pattabiraman. Formal security analysis of smart embedded systems. In *Proceedings of the 32nd Annual Conference on Computer Security Applications*, pages 1–15. ACM, 2016.
- [51] Carolyn Talcott. Cyber-physical systems and events. In Software-Intensive Systems and New Computing Paradigms, pages 101–115. Springer, 2008.
- [52] Uttam Thakore, Gabriel A Weaver, and William H Sanders. A quantitative methodology for security monitor deployment. In *Dependable* Systems and Networks (DSN), 2016 46th Annual IEEE/IFIP International Conference on, pages 1–12. IEEE, 2016.
- [53] John Villasenor. Next Homeland Security threat: cyber-physical attacks by drone strikes? https://www.brookings.edu/research/, 2011. [Online; accessed 19-July-2011].
- [54] Joachim Wegener and Matthias Grochtmann. Verifying timing constraints of real-time systems by means of evolutionary testing. *Real-Time Systems*, 15(3):275–298, 1998.
- [55] Westley Weimer and George C Necula. Mining temporal specifications for error detection. In International Conference on Tools and Algorithms for the Construction and Analysis of Systems, pages 461–476. Springer, 2005.
- [56] Peng Xie, Jason H Li, Xinming Ou, Peng Liu, and Renato Levy. Using bayesian networks for cyber security analysis. In *Dependable Systems* and Networks (DSN), 2010 IEEE/IFIP international conference on, pages 211–220. IEEE, 2010.

- [57] Ye Yan, Yi Qian, Hamid Sharif, and David Tipper. A survey on cyber security for smart grid communications. *IEEE Communications Surveys & Tutorials*, 14(4):998–1010, 2012.
- [58] Jinlin Yang, David Evans, Deepali Bhardwaj, Thirumalesh Bhat, and Manuvir Das. Perracotta: mining temporal api rules from imperfect traces. In *Proceedings of the 28th international conference on Software* engineering, pages 282–291. ACM, 2006.