

UNIVERSITY OF VERONA

Department of COMPUTER SCIENCE

Master's Degree in
COMPUTER SCIENCE AND ENGINEERING

Master Thesis

**Towards Process Comprehension of Industrial
Control Systems: a Framework for
Reverse-engineering Industrial Processes**

Candidate:

Marco OLIANI
VR457249

Supervisor:

Prof. Massimo MERRO

Co-supervisor:

Prof. Ruggero LANOTTE
University of Insubria

Academic Year 2022/2023

*“Someone cracked my password.
Now I need to rename my puppy.”*
(Unknown)

Abstract

The advent of Industry 4.0 has brought significant transformations to Industrial Control Systems (ICS), which monitor, coordinate, control, and integrate physical and engineered systems through computing and communication cores. The convergence of *Information Technology* (IT) and *Operational Technology* (OT) in ICS has improved operational efficiency but also exposed systems to cybersecurity vulnerabilities. Cyber-physical attacks targeting ICS aim to manipulate or disrupt their normal operation, posing risks to integrity and functionality.

Process comprehension plays a crucial role in executing successful cyber-physical attacks. Attackers must understand the operational domain, system architecture, industrial protocols, operational processes, process dependencies, attack surface, and attack goals. Techniques like probing, Man-in-the-Middle interception, and reverse engineering aid in gaining insights into system behavior without direct physical intervention.

This thesis analyzes Ceccato et al.'s reverse engineering-based methodology for ICS analysis [12], identifying limitations in its application to real-world scenarios. The framework's flexibility, network analysis, pre-processing, graphical analysis, invariant analysis, and business process analysis are improved to address these limitations and new functionalities are added. The enhanced framework is tested on a complex case study, the iTrust SWaT system, a scaled-down water treatment plant.

By understanding the process comprehension of ICS from an attacker's perspective and enhancing analysis methodologies, this research aims to contribute to the development of robust security measures in the context of Industry 4.0.

Contents

1	Introduction	1
2	Background on Industrial Control Systems	9
2.1	Industrial Control Systems Architecture	11
2.2	Operational Technology Networks	13
2.2.1	Field I/O Devices Layer	13
2.2.2	Controller Network Layer	14
2.2.2.1	Programmable Logic Controllers	14
2.2.2.2	Remote Terminal Units	18
2.2.3	Area Control Layer	19
2.2.3.1	Supervisory Control And Data Acquisition	19
2.2.3.2	Human-Machine Interface	19
2.2.4	Operations/Control Layer	20
2.2.5	Demilitarized Zone	20
2.2.6	Industrial Protocols	21
2.2.6.1	Modbus	22
2.2.6.2	EtherNet/IP	25
2.2.6.3	Common Industrial Protocol (CIP)	27
3	State of the Art	29
3.1	Literature on Process Comprehension	30
3.2	Ceccato et al.'s black-box dynamic analysis for water-tank systems	32

3.2.1	Testbed	34
3.2.2	Scanning of the System and Data Pre-processing . . .	37
3.2.3	Graphs and Statistical Analysis	39
3.2.4	Invariant Inference and Analysis	40
3.2.5	Business Process Mining and Analysis	42
3.2.6	Application	44
3.2.7	Limitations	50
4	Extending and Generalizing Ceccato et al.'s Framework	61
4.1	The Proposed New Framework	62
4.1.1	Framework Structure	65
4.1.2	Python Libraries and External Tools	66
4.2	Analysis Phases	67
4.2.1	A Little Testbed: Stage 1 of iTrust SWaT System . . .	68
4.2.2	Phase 0: Network Analysis	69
4.2.2.1	Extracting Data from PCAP Files	70
4.2.2.2	Network Information	72
4.2.3	Phase 1: Data Pre-processing	75
4.2.3.1	Subsystem Selection	76
4.2.3.2	Dataset Enrichment	79
4.2.3.3	Datasets Merging	86
4.2.3.4	Preliminary Analysis of the Obtained Sub-system	88
4.2.4	Phase 2: Graphs and Statistical Analysis	91
4.2.5	Phase 3: Invariant Inference and Analysis	94
4.2.5.1	Revised Daikon Output	95
4.2.5.2	Types of Analysis	99
4.2.6	Phase 4: Business Process Analysis	104
4.2.6.1	Process Mining of the Physical Process . . .	105
4.2.6.2	Network Data	109
4.2.7	Summary	110

5 Case study: the iTrust SWaT System	113
5.1 Architecture	113
5.1.1 Physical Process	114
5.1.2 Control and Communication Network	116
5.2 Datasets	118
5.2.1 Our Case Study: the 2015 Dataset	119
6 Our Framework at Work on the iTrust SWaT System	123
6.1 Preliminary Operations	124
6.2 Planning the Analysis Strategy	124
6.3 Reverse Engineering of the iTrust SWaT System	126
6.3.1 Reverse Engineering of PLC1 and PLC2	128
6.3.1.1 Pre-processing - Preliminary Analysis . . .	128
6.3.1.2 Graphs and Statistical Analysis	132
6.3.1.3 Invariant Inference and Analysis	135
6.3.1.4 Business Process Mining and Analysis . .	138
6.3.1.5 Properties	141
6.3.2 Reverse Engineering of PLC2 and PLC3	143
6.3.2.1 Pre-processing - Preliminary Analysis . .	143
6.3.2.2 Graphs and Statistical Analysis	147
6.3.2.3 Invariant Inference and Analysis	151
6.3.2.4 Business Process Mining and Analysis . .	154
6.3.2.5 Properties	157
6.3.3 Reverse Engineering of PLC3 and PLC4	158
6.3.3.1 Pre-processing - Preliminary Analysis . .	158
6.3.3.2 Graphs and Statistical Analysis	162
6.3.3.3 Invariant Inference and Analysis	164
6.3.3.4 Business Process Mining and Analysis . .	165
6.3.3.5 Properties	167
6.4 PLCs Architecture	168
7 Conclusion and Future Work	171
List of Figures	177

List of Tables	181
Listings	183
References	185

Introduction

¹ **T**HE advent of Industry 4.0 has sparked significant transformations in ² the realm of *Industrial Control Systems* (ICS), physical and engineered ³ systems whose operations are monitored, coordinated, controlled, and ⁴ integrated by a computing and communication core [1]. They are used ⁵ to control industrial processes such as manufacturing, product handling, ⁶ production, and distribution [2]. In recent years, there has been a rapid ex- ⁷ pansion in the interconnectivity and convergence of IT (*Information Tech-* ⁸ *nology*) and OT (*Operational Technology*) systems. While this integration ⁹ offers undeniable benefits in terms of operational efficiency and effective- ¹⁰ ness of ICSs, it has also exposed these systems to the prevalent vulnerabil- ¹¹ ities commonly associated with IT environments. Consequently, there has ¹² been a parallel rise in **cyber-physical attacks**, which originates in the cy- ¹³ ber environment and target the physical processes of these systems. These ¹⁴ attacks aim to manipulate or disrupt the normal operation of ICSs, posing ¹⁵ a considerable risk to their integrity and functionality.

¹⁶ **Contextualization** As mentioned earlier, the distinction between the IT ¹⁷ (*Information Technology*) and OT (*Operational Technology*) components within ¹⁸ an ICS is becoming increasingly blurred. The IT part pertains to the **cyber** ¹⁹ aspects of an industrial system, facilitating the collaboration and commu- ²⁰ nication of information processing technologies as well as technologies for

21 monitoring, control, and maintenance purposes [3].

22 On the other hand, the OT part primarily focuses on the aspects re-
23 lated to the **physical system**. The OT networks of ICSs typically encom-
24 pass devices, systems, networks, and controllers used for operating and
25 automating industrial processes. These networks consist of *field devices*,
26 such as sensors and actuators, which monitor and control the progression
27 of a physical process over time. They also include *Programmable Logic Con-*
28 *trollers* (PLCs) responsible for controlling the field devices, as well as one
29 or more *Human Machine Interfaces* (HMIs) that enable human operators to
30 interact with the PLCs and display status information and historical data
31 collected by the field devices within the ICS environment.

32 OT networks encompass two primary sub-networks: the *supervisory*
33 *control network*, which connects the PLCs and HMIs, and the *field commu-*
34 *nications network*, which establishes connections between the PLCs and the
35 associated field devices.

36 *Programmable Logic Controllers* (PLCs), which are essential elements of
37 ICSs, play a vital role in managing electrical equipment like pumps, valves,
38 centrifuges, and more. Serving as a bridge between the cyber and physi-
39 cal realms, PLCs possess a straightforward structure comprising a central
40 processing module (CPU) and additional modules for handling physical
41 inputs and outputs. The user program operates on the CPU and carries
42 out *scan cycles* that involve tasks such as gathering sensor data, execut-
43 ing controller code, and sending commands to actuator devices. Despite
44 modern controllers incorporating security measures to upload authorized
45 firmware, exploiting PLCs can still result in grave consequences for ICSs
46 including physical damage, operational disruptions, environmental haz-
47 ards, and even threats to human safety. These threats can originate from
48 various sources, such as hackers, insider threats, or state-sponsored actors,
49 and can have significant impacts on operational continuity, safety, and fi-
50 nancial losses. The potential impact of an ICS breach necessitates robust
51 security measures.

52 The increasing convergence of the IT and OT components in an ICS has

led to a transition from serial-based communication protocols to IP-based protocols. These protocols, known as **industrial protocols**, have become prevalent in modern industrial systems. Examples of industrial protocols include Modbus [4], DNP3 [5], EtherNet/IP [6], OPC UA [7], and S7comm [8] (which is a proprietary and closed protocol).

The adoption of IP-based protocols introduces a shared vulnerability between IT and OT networks. When OT networks utilize exposed networks like the Internet for communication, they become susceptible to the same security risks as IT networks. Similar to PLCs, breaches within the communication network can have significant consequences. Therefore, it is crucial to implement robust security measures to prevent potential cyber-physical attacks and ensure the integrity and security of the entire system.

Open Problems To execute a successful cyber-physical attack, an attacker must possess a comprehensive understanding of the target system's characteristics and behavior, commonly referred to as *process comprehension* [9]. This entails possessing knowledge about various aspects, including the **operational domain** (such as water distribution or power generation), the controllers used (such as PLCs, RTUs, etc.), the **network topology** associated with them, relevant **measurements** within the facility (such as pressure, temperature, etc.), as well as the exposed physical components like **sensors and actuators** that can be vulnerable to attacks. By attaining such insights, the attacker can effectively identify vulnerabilities, exploit weaknesses, and manipulate the system in a targeted and *stealthy* manner. Moreover, to bypass certain types of *intrusion detection systems* (IDSs), the attacker needs to possess a general understanding of the **physical invariants** of the system.

Understanding the process operations and functionalities of ICS from an attacker's perspective is a crucial aspect of conducting targeted cyber-attacks. By comprehending how the systems operate, their vulnerabilities, and potential impact, attackers can devise effective strategies to infiltrate

84 and compromise the ICS environment. Here are some key points related
85 to process comprehension of ICS from an attacker's viewpoint:

- 86 • **System Architecture:** attackers aim to understand the architecture
87 of the ICS, including the components involved, such as sensors, actuators,
88 controllers, and human-machine interfaces. This knowledge
89 helps them identify potential entry points and vulnerabilities within
90 the system;
- 91 • **Industrial Protocols:** familiarity with industrial communication pro-
92 tocols, such as Modbus [4], DNP3 [5], or EtherNet/IP - CIP [10], is
93 crucial for attackers. Understanding the protocols enables them to
94 analyze the communication between different ICS components and
95 potentially exploit vulnerabilities or manipulate the data flow;
- 96 • **Operational Processes:** attackers seek to comprehend the operational
97 processes and workflows within the targeted ICS. This includes un-
98 derstanding the sequence of actions, system states, and interactions
99 between various components. Such knowledge enables attackers to
100 identify critical points where disruptions or malicious actions can
101 have significant consequences;
- 102 • **Process Dependencies:** attackers analyze the dependencies between
103 different physical processes or systems within the ICS. They aim to
104 identify dependencies that, if disrupted or compromised, could have
105 a cascading effect on the overall operation. This helps them strate-
106 gize targeted attacks for maximum impact;
- 107 • **Attack Surface:** by evaluating the attack surface of the ICS, includ-
108 ing network connectivity, remote access options, and third-party in-
109 tegrations, attackers can identify potential entry points and avenues
110 for exploitation;
- 111 • **Attack Goals:** attackers define their specific goals, whether it be dis-
112 ruption of operations, theft of sensitive data, sabotage, or any other
113 malicious intent. Understanding the desired outcomes helps them

114 tailor their attack strategies and prioritize their actions within the
115 ICS environment.

116 Attackers can acquire a good level of process comprehension of an
117 ICS through several methods. For instance, they can utilize **probing** tech-
118 niques [11], which involve introducing slight perturbations to the system
119 and observing its response and subsequent return to its original state. By
120 analyzing these reactions, attackers can gain valuable insights into the sys-
121 tem's behavior.

122 Another technique is the **Man-in-the-Middle approach** [9], where at-
123 tackers intercept and analyze network communications to gather infor-
124 mation about the system. This allows them to conduct **reconnaissance** of
125 the system while avoiding detection.

126 **Reverse engineering** [12][13] is another method attackers can employ.
127 By scanning memory logs of network-exposed PLCs and analyzing net-
128 work traffic related to the industrial protocols, they can approximate the
129 physical system's model. This technique enables them to derive useful in-
130 formation without directly intervening in the physical system, as they can
131 work offline using collected logs and analysis tools.

132 By leveraging these techniques, attackers can execute targeted and pre-
133 cise attacks on the system, avoiding the need for *Denial of Service* (DoS)
134 tactics. This enables them to focus on exploiting specific vulnerabilities
135 and achieving their objectives with greater precision and effectiveness.

136 Contribution

137 The primary objective of this thesis is to conduct an analysis of Ceccato
138 et al.'s reverse engineering-based methodology [12]. The aim is to examine
139 its strengths and limitations and apply it to a real-world case study.

140 Ceccato et al.'s methodology relies on a specialized framework that
141 performs the analysis of an ICS through multiple phases. It begins with
142 the collection of system-related data by scanning the logs of exposed PLCs

for information on the physical system and analyzing network traffic related to the industrial protocols used for communications. Subsequent phases involve enriching the obtained data for the physical system, conducting graphical-statistical analysis of PLC logs, inferring system invariants using Daikon, and understanding the overall behavior of the PLCs in conjunction with network traffic data through a business process analysis. Ceccato et al. utilize a virtualized testbed, simulating a water treatment plant, with communication occurring through the Modbus protocol.

Ceccato et al.'s framework has **limitations** that hinder its use on systems other than their specific testbed. One limitation is that the testbed itself is overly simplified compared to real-world scenarios, with only three stages and few components, resulting in a limited number of registers. Moreover, the virtualization does not account for the real physics of water and its oscillations. Another limitation is that their Graphical Analysis can display only one register at a time, making it impossible to have an overall view of the system or capture relationships between registers. Additionally, the output produced by the external tool Daikon [14] for invariant analysis is difficult to interpret due to its poorly readable format.

Overall, the Ceccato et al. framework lacks flexibility in analyzing different elements and relies heavily on ad hoc solutions for their specific testbed and the Modbus protocol. This can lead to limited and inaccurate results when applied to other industrial systems, and in some cases, it may even be impossible to perform the analysis.

The objective of this thesis is to address these limitations identified in Ceccato et al.'s framework by improving and expanding upon the original framework. This involves a complete overhaul and rewriting of the code to introduce previously absent features and enhance existing ones. The framework has been made independent of the system being analyzed, allowing for customization through a configuration file. Furthermore, enhancements and expansions have been introduced across all analysis phases.

For instance, *Network Analysis* has been added to discover **network**

174 **topology and device communications**, regardless of the industrial pro-
175 tocol used. Enhancements have also been made to other phases, including
176 improvements to the **command-line interface** for increased flexibility and
177 user-friendliness, as well as the ability to extend the scope of the analysis.

178 The *Pre-processing* phase, responsible for enriching data obtained from
179 physical system scans, now allows the association of additional fields such
180 as slope with specific registers or groups of registers. This provides a more
181 accurate dataset and facilitates operations in other phases. Additionally,
182 a new **preliminary analysis phase** has been introduced, automating the
183 recognition of probable measurements and actuators, and analyzing vari-
184 ous properties related to them.

185 The *Graphical Analysis* component has undergone significant improve-
186 ments to enhance its effectiveness. **Subplots** now enable the visualization
187 of registers as a whole, facilitating the identification of relationships be-
188 tween them. Users can select specific registers to view and zoom in on
189 specific areas of the graphs without losing information.

190 *Invariant Analysis* has also been revised to **improve the output gener-
191 ated by Daikon**, making it more concise and readable for easier identifi-
192 cation of invariants. **Two semi-automatic analyses** based on individual
193 actuator states and current system states have been introduced as well.

194 Finally, the *Business Process Analysis* has been revamped to extract as
195 much information as possible from the physical system, including actual
196 system states. This information is then integrated with network data using
197 a newly developed solution instead of relying on an external tool as in the
198 previous version.

199 The new framework will be tested on a **more complex case study**,
200 namely the iTrust SWaT system [15], which is an actual water treatment
201 system. Compared to Ceccato et al.'s testbed, the iTrust SWaT system in-
202 volves more PLCs, registers, and components, presenting additional chal-
203 lenges. For instance, fluctuations in water levels within the tanks may
204 introduce data perturbations that need to be mitigated to ensure accurate
205 results aligned with reality.

206 Furthermore, the iTrust SWaT system employs the CIP over EtherNet/IP
207 protocol, which differs significantly from Modbus. Consequently, differ-
208 ent implementation approaches are required when considering network
209 data within the various analysis steps, as compared to Ceccato et al.'s
210 methodology.

211 **Outline**

212 The thesis is structured as follows:

213 **Chapter 2:** provides a background on Industrial Control Systems, (ICSS)
214 describing their structure, components, and some of the network
215 communication protocols used;

216 **Chapter 3:** after an introductory section that provides a brief overview of
217 the existing literature on process comprehension in industrial con-
218 trol systems, this chapter focuses on a specific paper that outlines
219 a methodology to attaining process comprehension of an industrial
220 system by employing dynamic blackbox analysis;

221 **Chapter 4:** defines a proposal to improve and extend the methodology
222 outlined in the previous chapter;

223 **Chapter 5:** presents the case study on which the proposed methodology
224 will be applied;

225 **Chapter 6:** shows how the proposed methodology is applied to the case
226 study illustrated above;

227 **Chapter 7:** outlines final conclusions and future work.

Background on Industrial Control Systems

²²⁸ INDUSTRIAL control systems (ICSs) are physical and engineered systems
²²⁹ whose operations are monitored, coordinated, controlled and integrated
²³⁰ by a computing and communication core [1]. They are used to control in-
²³¹ dustrial processes such as manufacturing, product handling, production,
²³² and distribution [2]. ICSs are often found in critical infrastructure facilities
²³³ such as power plants, oil and gas refineries, and chemical plant.

²³⁴ ICSs are different from traditional IT systems in several key ways. Firstly,
²³⁵ ICSs are designed to control physical processes, whereas IT systems are
²³⁶ designed to process and store data. This means that ICSs have different
²³⁷ requirements for availability, reliability, and performance. Secondly, ICSs
²³⁸ are typically deployed in environments that are harsh and have limited
²³⁹ resources, such as extreme temperatures and limited power. Thirdly, the
²⁴⁰ protocols hardware and software used in ICSs are often proprietary.

²⁴¹ ICSs are becoming increasingly connected to the internet and other net-
²⁴² works, which has led to increased concerns about their security. Industrial
²⁴³ systems were not originally designed with security in mind, and many of
²⁴⁴ them have known vulnerabilities that could be exploited by attackers. Ad-
²⁴⁵ ditionally, the use of legacy systems and equipment can make it difficult
²⁴⁶ to implement security measures. As a result, ICSs are increasingly seen
²⁴⁷ as a potential target for cyber attacks, which could have serious conse-
²⁴⁸ quences for the safe and reliable operation of critical infrastructure: some

notorious examples of cyber attacks are (i) the **STUXnet** worm [16], which purpose was to sabotage the nuclear centrifuges of the enrichment plant at the Natanz nuclear facility in Iran; (ii) **Industroyer** [17], also referred as *Crashoverride*, responsible for the attack on the Ukrainian power grid on December 17, 2016; (iii) the attack on February, 2021 to a water treatment plant in Oldsmar, Florida [18], where the level of sodium hydroxide was intentionally increased to a level approximately 100 times higher than normal.

The increasing connectivity of ICSs and the associated security risks have led to a growing interest in the field of ICS security. Researchers and practitioners are working to develop new security technologies, standards, and best practices to protect ICSs from cyber attacks. This includes efforts to improve the security of ICS networks and devices, as well as the development of new monitoring and detection techniques to identify and respond to cyber attacks.

Table 2.1 summarizes the differences between traditional IT and ICSs [3]:

	Traditional IT	ICSs
Focus	Data	Asset
Update Frequency	High	Low
	Confidentiality	Availability
Priority	Integrity	Integrity
	Availability	Confidentiality
Operating System	Standardized	Proprietary
Protocols	Standardized	Proprietary
Attacker Motivation	Monetization	Disruption

Table 2.1: Differences between Information Technology (IT) and Industrial Control Systems (ICSs)

2.1 Industrial Control Systems Architecture

In the past, there has been a clear division between *Information Technology* (IT) and *Operational Technology* (OT), both at the technical and organizational levels. Each domain has maintained its own distinct technology stacks, protocols, and standards. However, with the emergence of Industry 4.0 and the rapid expansion of industrial automation, which heavily relies on IT tools for monitoring and controlling critical infrastructures, the boundary between IT and OT has started to blur. This trend has paved the way for greater integration between these two domains, thus improving productivity and process quality.

General ICS architecture consists in **six levels** each representing a functionality: this architecture comprising the OT and IT parts is represented in Figure 2.1 [19][3], according to the *Purdue Enterprise Reference Architecture* (PERA), or simply **Purdue Model**:

- Level 0 (**Processes, or Field I/O Devices**): contains **field devices**.
- Level 1 (**Intelligent Devices, or Controller Network**): includes **local or remote controllers** that sense, monitor and control the physical process, such as **PLCs** (2.2.2.1) and **RTUs** (2.2.2.2). Controllers interface directly to the field devices reading data from sensors and sending commands to actuators.
- Level 2 (**Control Systems, or Area Control**): contains computer systems used to supervising and monitoring the physical process: they provide a **Human-Machine Interface** (*HMI*, 2.2.3.2) and *Engineering Workstations* (EW) for operator control.
- Level 3 (**Manufacturing/Site Operations, or Operations/Control**): comprises systems used to manage the production workflow for plant-wide control: they collate informations from the previous levels and store them in Data Historian servers.

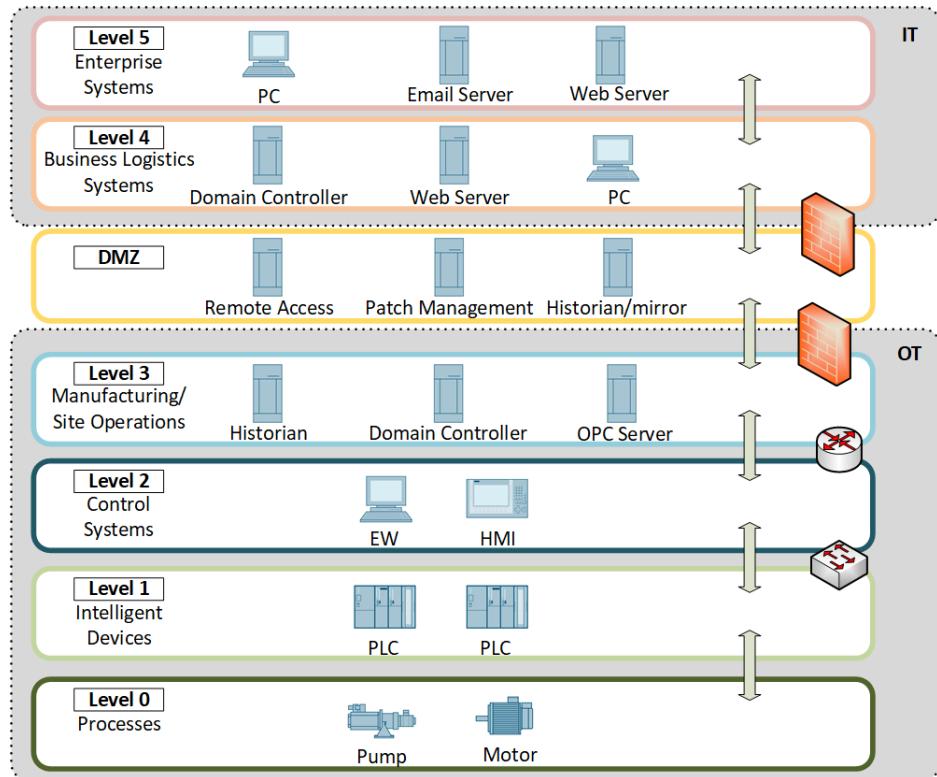


Figure 2.1: ICS architecture schema

- 293 • **Industrial Demilitarized Zone (DMZ):** intermediate level that connects the *Operational Technology* (OT) part (levels 0-3) with the *Information Technology* (IT) part of the system (levels 4 and 5). Communication takes place indirectly through services such as *proxy servers* and *remote access servers*, which act as intermediaries between the two environments.
- 294
- 295
- 296
- 297
- 298
- 299 • Level 4 (**Business Logistics Systems, or Business Planning/Logistics**): collect and aggregates data from the Manufacturing/Site Operations level overseeing the IT-related activities to generate **reporting** to the Enterprise System layer. At this layer we can find application and e-mail servers, and *Enterprise Resource Planning* (ERP) systems.
- 300
- 301
- 302
- 303
- 304 • Level 5 (**Enterprise Systems**): represents the enterprise network, used for the business-to-business activities and for business-to-client pur-
- 305

306 pose services. At Enterprise Systems level are typical IT services such
307 as mail servers, web servers and all the systems used to manage the
308 ongoing process.

309 As previously discussed, the gap between IT and OT is steadily nar-
310 rowing. Nowadays, it is increasingly common to encounter IT elements
311 within the OT realm. For example, desktop PCs are now frequently found
312 in OT environments, and industrial devices are interconnected using stan-
313 dard IT communication protocols like TCP and UDP.

314 2.2 Operational Technology Networks

315 *Operational Technology* primarily encompasses the **tangible aspects** of
316 Industrial Control Systems and directly interfaces with the physical pro-
317 cesses of the monitored systems. Its main purpose is to **manage and con-**
318 **trol the procedures** involved in creating and correcting physical value in
319 various equipment.

320 This section will focus on the key aspects and components of Opera-
321 tional Technology network, with specific reference to the first four levels
322 of the Purdue model previously seen.

323 2.2.1 Field I/O Devices Layer

324 This level concerns all aspects related to the physical environment and
325 the physical elements that are part of it, which have the ability to actively
326 influence the environment.

327 These physical elements are represented by **Field Devices**, i.e., **sensors**
328 and **actuators** used to collect data from the process and control it: sen-
329 sors are the elements responsible for reading specific values related to the
330 physical environment (e.g., the level of a liquid), while actuators change
331 its behavior and characteristics (e.g., opening or closing a valve to make
332 the liquid flow). Examples of field devices include temperature sensors,
333 pressure sensors, valves and pumps.

334 2.2.2 Controller Network Layer

335 *Controller Network* layer includes devices that handle data from and
336 to the *Field I/O Devices* layer. This kind of device is capable of gathering
337 data from sensors, updating its internal state, and activating actuators (for
338 example opening or close a pump that controls the level of a tank), making
339 decisions based on a customized program, known as its control logic.

340 Commonly found within this layer are *Programmable Logic Controllers*
341 (PLCs) and *Remote Terminal Units* (RTUs): in the upcoming sections, we
342 will examine these elements in detail.

343 2.2.2.1 Programmable Logic Controllers

344 A *Programmable Logic Controller* (PLC) is a **small and specialized in-**
345 **dustrial computer** having the capability of controlling complex industrial
346 and manufacturing processes [20].

347 Compared to relay systems and personal computers, PLCs are opti-
348 mized for control tasks and industrial environments: they are rugged and
349 designed to withdraw harsh conditions such as dust, vibrations, humid-
350 ity and temperature: they have more reliability than personal computers,
351 which are more prone to crash, and they are more compact a require less
352 maintenance than a relay system.

353 Furthermore, I/O interfaces are already on the controller, so PLCs are
354 easier to expand with additional I/O modules (if in a rack format) to man-
355 age more inputs and ouputs, without reconfiguring hardware as in relay
356 systems when a reconfiguration occurs.

357 PLCs are more *user-friendly*: they are not intended (only) for computer
358 programmers, but designed for engineers with a limited knowledge in
359 programming languages: control program can be entered with a simple
360 and intuitive language based on logic and switching operations instead of
361 a general-purpose programming language (*i.e.* C, C++, ...).

- 362 **PLC Architecture** The basic hardware architecture of a PLC consists of
363 these elements [21]:

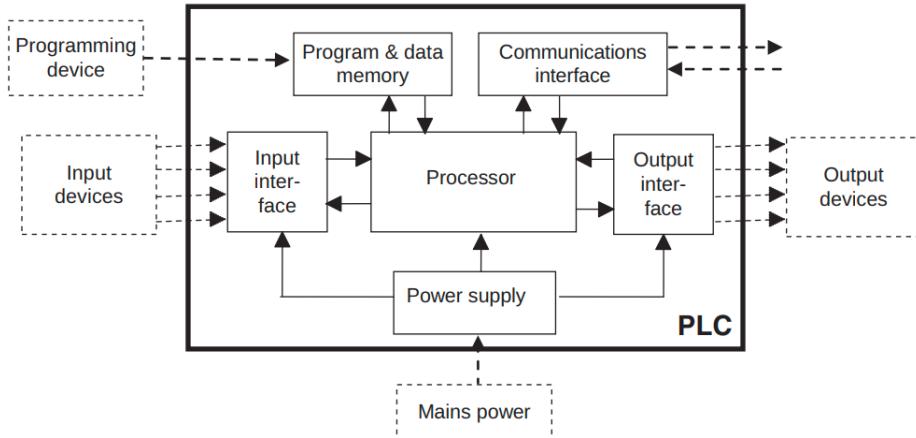


Figure 2.2: PLC architecture

- **Processor unit (CPU):** contains the microprocessor. This unit interprets the input signals from I/O modules, executes the control program stored in the Memory Unit and sends the output signals to the I/O Modules. The processor unit also sends data to the Communication interface, for the communication with additional devices.
- **Power supply unit:** converts AC voltage to low DC voltage.
- **Programming device:** is used to store the required program into the memory unit.
- **Memory Unit:** consists in RAM memory and ROM memory. RAM memory is used for storing data from inputs, ROM memory for storing operating system, firmware and user program to be executed by the CPU.
- **I/O modules:** provide interface between sensors and final control elements (actuators).
- **Communications interface:** used to send and receive data on a network from/to other PLCs.

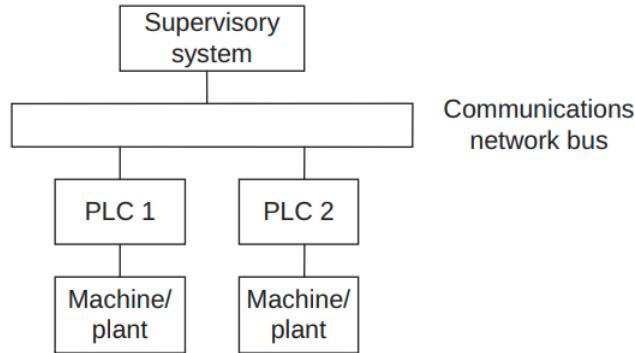


Figure 2.3: PLC communication schema

380 **PLC Programming** Two different programs are executed in a PLC: the
381 **operating system** and the **user program**.

382 The operating system tasks include executing the user program, man-
383 aging memory areas and the *process image table* (memory registers where
384 inputs from sensors and outputs for actuators are stored).

385 The user program needs to be uploaded on the PLC via the program-
386 ming device and runs on the process image table in *scan cycles*: each scan
387 is made up of three phases [12]:

- 388 1. reading inputs from the process images table
- 389 2. execution of the control code and computing the physical process evolution
- 390 3. writing output to the process image table to have an effect on the physical process. At the end of the cycle, the process image table is refreshed by the CPU

394 Standard PLCs **programming languages** are basically of two types:
395 **textuals** and **graphicals**. Textual languages include languages such as
396 *Instruction List* (IL) and *Structured Text* (ST), while *Ladder Diagrams* (LD),
397 *Function Block Diagram* (FBD) and *Sequential Function Chart* (SFC) belong
398 to the graphical languages.

399 Graphical languages are more simple and immediate comparing to the
 400 textual ones and are preferred by programmers because of their features
 401 and simplicity, in particular the **Ladder Logic programming** (see Figure
 402 2.4 for a comparison).

```

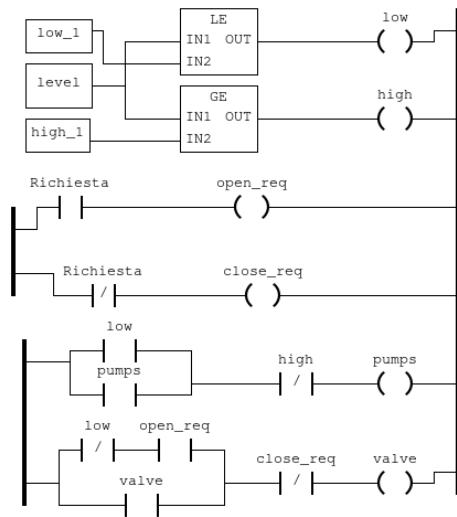
PROGRAM PLC1
VAR
  level AT %IW0 : INT;
  Richiesta AT %QX0..2 : BOOL;
  request AT %IW1 : INT;
  pumps AT %QX0..0 : BOOL;
  valve AT %QX0..1 : BOOL;
  low AT %MW0..0 : BOOL;
  high AT %MW0..1 : BOOL;
  open_req AT %MW0..3 : BOOL;
  close_req AT %MW0..4 : BOOL;
  low_1 AT %MW0 : INT := 40;
  high_1 AT %MW1 : INT := 80;
END_VAR
VAR
  LE3_OUT : BOOL;
  GE7_OUT : BOOL;
END_VAR

LE3_OUT := LE(level, low_1);
low := LE3_OUT;
GE7_OUT := GE(level, high_1);
high := GE7_OUT;
open_req := Richiesta;
close_req := NOT(Richiesta);
pumps := NOT(high) AND (low OR pumps);
valve := NOT(close_req) AND (open_req AND NOT(low) OR valve);

END_PROGRAM

CONFIGURATION Config0
  RESOURCE Res0 ON PLC
    TASK task0(INTERVAL := T#20ms,PRIORITY := 0);
    PROGRAM instance0 WITH task0 : PLC1;
  END_RESOURCE
END_CONFIGURATION
  
```

(a) Example of ST programming



(b) Example of Ladder Logic

Figure 2.4: Comparison between ST language and Ladder Logic

403 **PLC Security** PLCs were originally designed to operate as closed sys-
 404 tems, not connected and exposed to the outside world via communication
 405 networks: the question of the safety of these systems, therefore, was not
 406 a primary aspect. The advent of Internet has brought undoubted advan-
 407 tages, but has introduced problems relating to the safety and protection of
 408 PLCs from external attacks and vulnerabilities.

409 Indeed, a variety of different communication protocols used in ICSs are
 410 designed to be efficient in communications, but do not provide any secu-
 411 rity measure i.e. confidentiality, authentication and data integrity, which
 412 makes these protocols vulnerable against many of the IT classic attacks
 413 such as *Replay Attack* or *Man in the Middle Attack*.

414 Countermeasures to enhance security in PLC systems may include [22]:

- 415 • protocol modifications implementing **data integrity, authentication**
416 and **protection** against *Replay Attacks*
- 417 • use of *Intrusion Detection and Prevention Systems* (IDP)
- 418 • creation of *Demilitarized Zones* (DMZ) on the network

419 In addition to this, keeping the process network and Internet sepa-
420 rated, limiting the use of USB devices among users to reduce the risks of
421 infections, and using strong account management and maintenance poli-
422 cies are best practices to prevent attacks and threats and to avoid potential
423 damages.

424 2.2.2.2 Remote Terminal Units

425 *Remote Terminal Units* (RTUs) are computers with radio interfacing sim-
426 ilar to PLCs: they transmit telemetry data to the control center or to the
427 PLCs and use messages from the master supervisory system to control
428 connected objects [23].

429 The purpose of RTUs is to operate efficiently in remote and isolated
430 locations by utilizing wireless connections. In contrast, PLCs are designed
431 for local use and rely on high-speed wired connections. This key difference
432 allows RTUs to conserve energy by operating in low-power mode for ex-
433 tended periods using batteries or solar panels. As a result, RTUs consume
434 less energy than PLCs, making them a more sustainable and cost-effective
435 option for remote operations.

436 Industries that require RTUs often operate in areas without reliable ac-
437 cess to the power grid or require monitoring and control substations in re-
438 mote locations. These include telecommunications, railways, and utilities
439 that manage critical infrastructure such as power grids, pipelines, and wa-
440 ter treatment facilities. The advanced technology of RTUs allows these in-
441 dustries to maintain essential services, even in challenging environments
442 or under adverse weather conditions.

2.2.3 Area Control Layer

The Area Control layer encompasses hardware and software systems useful for supervising, monitoring and controlling the physical process, driving the behavior of the entire infrastructure. The layer includes systems such as *Supervisory Control and Data Acquisition* (SCADA), *Distributed Control Systems* (DCSs), that perform SCADA functions but are usually deployed locally, and engineer workstations.

2.2.3.1 Supervisory Control And Data Acquisition

Supervisory Control And Data Acquisition (SCADA) is a system of software and hardware elements that allows industrial organizations to [24]:

- Control industrial processes locally or at remote locations;
- Monitor, gather, and process real-time data;
- Directly interact with devices such as sensors, valves, pumps, motors, and more through human-machine interface (HMI) software;
- Record and aggregate events to send to historian server.

The SCADA software processes, distributes, and displays the data, helping operators and other employees analyze the data and make important decisions.

2.2.3.2 Human-Machine Interface

The *Human-Machine Interface* (HMI) is the hardware and software interface that operators use to monitor the processes and interact with the ICS. A HMI shows the operator and authorized users information about system status and history; it also allows them to configure parameters on the ICS such as set points and, send commands and make control decisions [25].

The HMI can be in the form of a physical panel, with buttons and indicator lights, or PC software as shown in Figure 2.5.

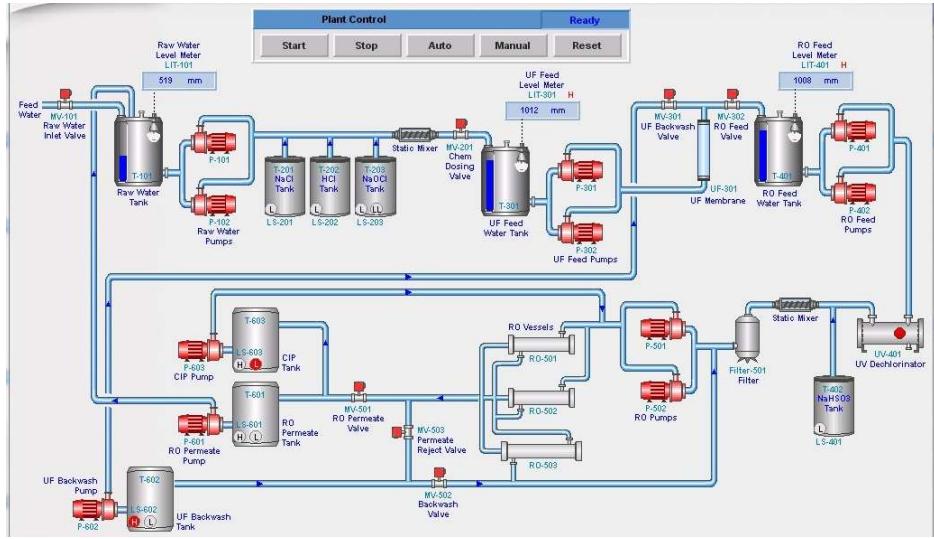


Figure 2.5: Example of HMI for a water treatment plant

2.2.4 Operations/Control Layer

Within this zone, there are specialized OT devices that are utilized to manage production workflows on the shop floor [26]. These devices include:

- Manufacturing Operations Management (MOM)* systems, which are responsible for overseeing production operations.
- Manufacturing Execution Systems (MES)*, which collect real-time data to optimize production processes.
- Data Historians*, which store process data and, in modern solutions, analyze it within its contextual framework.

2.2.5 Demilitarized Zone

This zone comprises security systems like firewalls, proxies, *Intrusion Detection and Prevention systems* (IDP) and *Security Information and Event Management* (SIEM) systems which are implemented to mitigate the risk of lateral threat movement between IT and OT domains. With the rise

485 of automation, the need for bidirectional data flows between OT and IT
486 systems has increased. The convergence of IT and OT in this layer can offer
487 organizations a competitive edge. However, it's important to note that
488 adopting a flat network approach in this context can potentially heighten
489 cyber risks for the organization.

490 **2.2.6 Industrial Protocols**

491 *Industrial Protocols* are the networks that are used to connect the dif-
492 ferent components of the ICS and allow them to communicate with each
493 other. Industrial Protocols can include wired and wireless networks, such
494 as Ethernet/IP, Modbus, DNP3, Profinet and others.

495 As mentioned at the beginning of this Chapter, industrial systems dif-
496 fer from classical IT systems in the purpose for which they are designed:
497 controlling physical processes the former, processing and storing data the
498 latter. For this reason, ICSs require different communication protocols
499 than traditional IT systems for real time communications and data trans-
500 fer.

501 A wide variety of industrial protocols exists: this is because originally
502 each vendor developed and used its own proprietary protocol. How-
503 ever, these protocols were often incompatible with each other, resulting
504 in devices from different vendors being unable to communicate with each
505 other.

506 To solve this problem, standards were defined with a view to allowing
507 these otherwise incompatible device to intercommunicates.

508 Among all the various protocols, some have risen to prominence as
509 widely accepted standards. These *de facto* protocols are commonly utilized
510 in industrial systems due to their proven reliability and effectiveness. In
511 the following sections, we will provide a brief overview of some of the
512 most prevalent and widely used protocols in the industry.

513 **2.2.6.1 Modbus**

514 *Modbus* is a serial communication protocol developed by Modicon (now
 515 Schneider Electric) in 1979 for use with its PLCs [4] and designed expressly
 516 for industrial use: it facilitates interoperability of different devices con-
 517 nected to the same network (sensors, PLCs, HMIs, ...) and it is also often
 518 used to connect RTUs to SCADA acquisition systems.

519 Modbus is the most widely used communication protocol among in-
 520 dustrial systems because it has several advantages:

- 521 • simplicity of implementation and debugging
- 522 • it moves raw bits and words, letting the individual vendor to repre-
 523 sent the data as it prefers
- 524 • it is, nowadays, an **open** and *royalty-free* protocol: there is no need
 525 to sustain licensing costs for implementation and use by industrial
 526 device vendors

527 Modbus is a **request/response** (or *master/slave*) protocol: this makes it
 528 independent of the transport layer used.

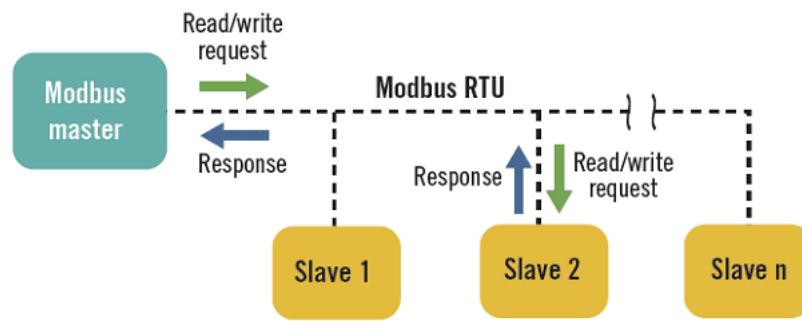


Figure 2.6: Modbus Request/Response schema

529 In this kind of architecture, a single device (master) can send requests
 530 to other devices (slaves), either individually or in broadcast: these slave
 531 devices (usually peripherals such as actuators) will respond to the master

532 by providing data or performing the action requested by the master using
 533 the Modbus protocol. Slave devices cannot generate requests to the master
 534 [27].

535 There are several variants of Modbus, of which the most popular and
 536 widely used are Modbus RTU (used in serial port connections) and Mod-
 537 bus TCP (which instead uses TCP/IP as the transport layer). Modbus TCP
 538 embeds a standard Modbus frame in a TCP frame (see Figure 2.7): both
 539 masters and slaves listen and receive data via TCP port 502.

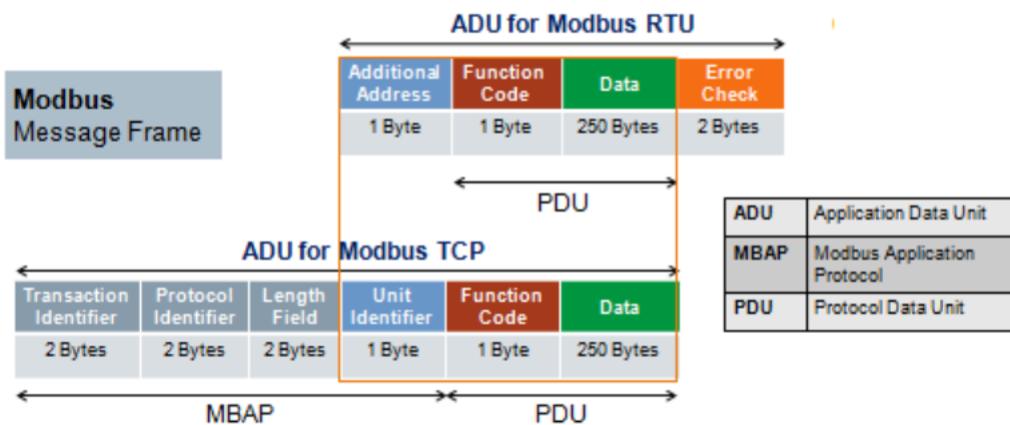


Figure 2.7: Modbus RTU frame and Modbus TCP frame

540 **Modbus registers** Modbus provides four object types, which map the
 541 data accessed by master and slave to the PLC memory:

- 542 • *Coil*: binary type, read/write accessible by both masters and slaves
- 543 • *Discrete Input*: binary type, accessible in read-only mode by masters
 544 and in read/write mode by slaves
- 545 • *Analog Input*: 16 bits in size (word), are accessible in read-only mode
 546 by masters and in read/write mode by slaves
- 547 • *Holding Register*: 16 bits in size (word), accessible in read/write mode
 548 by both masters and slaves. Holding Registers are the most com-
 549 monly used registers for output and as general memory registers.

550 **Modbus Function Codes** *Modbus Function Codes* are specific codes used
551 by the Modbus master within a request frame (see Figure 2.7) to tell the
552 Modbus slave device which register type to access and which action to
553 perform on it.

554 Two types of Function Codes exists: for data access and for diagnostic
555 Function Codes list for data access are listed in Table 2.2:

Function Code	Description
FC01	Read Coils
FC02	Read Discrete Input
FC03	Read Holding Registers
FC04	Read Analog Input Registers
FC05	Write/Force Single Coil
FC06	Write/Force Single Holding Register
FC15	Write/Force Multiple Coils
FC16	Write/Force Multiple Holding Registers

Table 2.2: Modbus Function Codes list

556 **Modbus Security Issues** Despite its simplicity and widespread use, the
557 Modbus protocol does not have any security feature, which exposes it to
558 vulnerabilities and attacks.

559 Data in Modbus are transmitted unencrypted (*lack of confidentiality*), with
560 no data integrity controls (*lack of integrity*) and authentication checks (*lack*
561 *of authentication*), in addition to the *lack of session*. Hence, the protocol is
562 vulnerable to a variety of attacks, such as Denial of Services (DoS), buffer
563 overflows and reconnaissance activities.

564 The easiest attack to bring to the Modbus protocol, however, is **packet**
565 **sniffing** (Figure 2.8): since, as mentioned earlier, network traffic is un-
566 encrypted and the data transmitted is in cleartext, it is sufficient to use
567 a packet sniffer to capture the network traffic, read the packets and thus

- 568 gather informations about the system such as ip addresses, function codes
 569 of requests and to modify the operation of the devices.

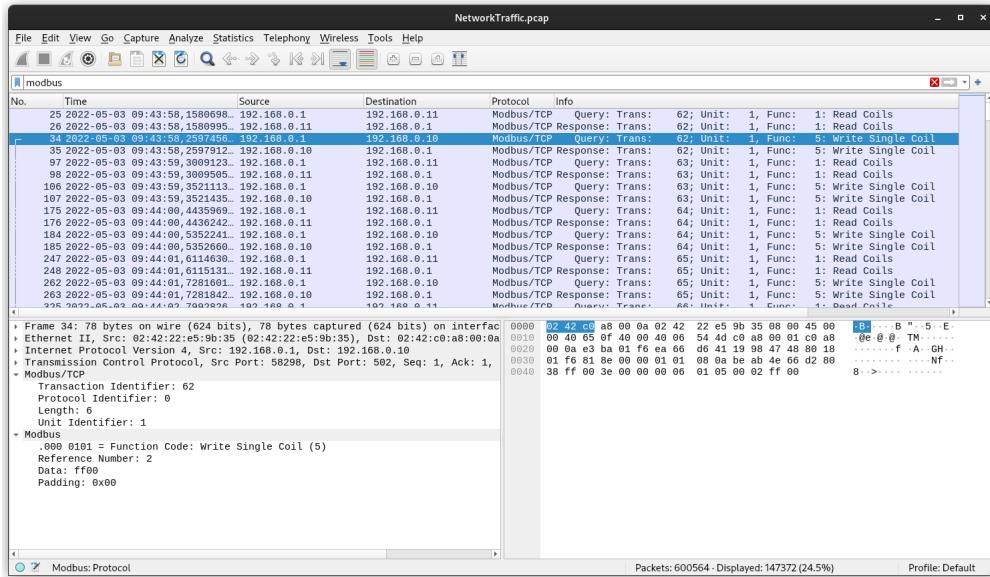


Figure 2.8: Example of packet sniffing on the Modbus protocol

570 To make the Modbus protocol more secure, an encapsulated version
 571 was developed within the *Transport Security Layer* (TLS) cryptographic
 572 protocol, also using mutual authentication. This version of the Modbus
 573 protocol is called **Secure Modbus** or **Modbus TLS**. In addition to this,
 574 Secure Modbus also includes X.509-type certificates to define permissions
 575 and authorisations [28].

576 2.2.6.2 EtherNet/IP

577 *EtherNet/IP* (where IP stands for *Industrial Protocol*) is an open indus-
 578 trial protocol that allows the *Common Industrial Protocol* (CIP) to run on a
 579 typical Ethernet network [10]. It is supported by ODVA [29].

580 EtherNet/IP uses the major Ethernet standards, such as IEEE 802.3 and
 581 the TCP/IP suite, and implements the CIP protocol stack at the upper lay-
 582 ers of the OSI stack (see Figure 2.9). It is furthermore compatible with the
 583 main Internet standard protocols, such as SNMP, HTTP, FTP and DHCP,

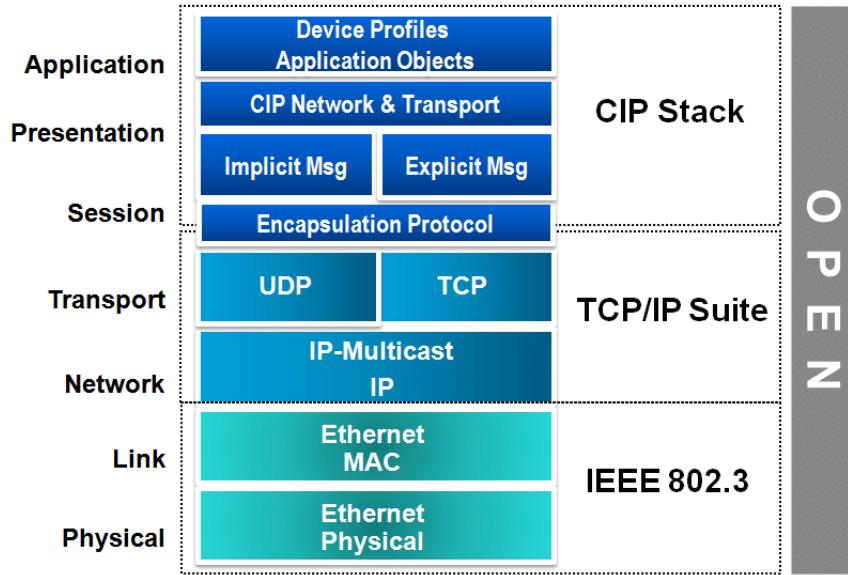


Figure 2.9: OSI model for EtherNet/IP stack

584 and other industrial protocols for data access and exchange such as *Open
585 Platform Communication* (OPC).

586 **Physical and Data Link layer** The use of the IEEE 802.3 standard allows
587 EtherNet/IP to flexibly adopt different network topologies (star, linear,
588 ring, etc.) over different connections (copper, fibre optic, wireless, etc.), as
589 well as the possibility to choose the speed of network devices.

590 IEEE 802.3 in addition defines at Data Link layer the *Carrier Sense Multiple
591 Access - Collision Detection* (CSMA/CD) protocol, which controls access
592 to the communication channel and prevents collisions.

593 **Transport layer** At the transport level, EtherNet/IP encapsulates mes-
594 sages from the CIP stack into an Ethernet message, so that messages can
595 be transmitted from one node to another on the network using the TCP/IP
596 protocol. EtherNet/IP uses two forms of messaging, as defined by CIP
597 standard [10][6]:

- 598 • **unconnected messaging:** used during the connection establishment
599 phase and for infrequent, low priority, explicit messages. Uncon-

600 nected messaging uses TCP/IP to transmit messages across the net-
601 work asking for connection resource each time from the *Unconnected*
602 *Message Manager* (UCMM).

- 603 • **connected messaging:** used for frequent message transactions or for
604 real-time I/O data transfers. Connection resources are reserved and
605 configured using communications services available via the UCMM.

606 EtherNet/IP has two types of message connection [10]:

- 607 – **explicit messaging:** *point-to-point* connections to facilitate *request-*
608 *response* transactions between two nodes. These connections use
609 TCP/IP service on port 44818 to transmit messages over Ether-
610 net.
- 611 – **implicit messaging:** this kind of connection moves application-
612 specific **real-time I/O data** at regular intervals. It uses multicast
613 *producer-consumer* model in contrast to the traditional *source-*
614 *destination* model and UDP/IP service (which has lower proto-
615 col overhead and smaller packet size than TCP/IP) on port 2222
616 to transfer data over Ethernet.

617 **Session, Presentation and Application layer** At the upper layers, Ether-
618 Net/IP implements the CIP protocol stack. We will discuss this protocol
619 more in detail in Section 2.2.6.3.

620

2.2.6.3 Common Industrial Protocol (CIP)

621 The *Common Industrial Protocol* (CIP) is an open industrial automation
622 protocol supported by ODVA. It is a **media independent** (or *transport in-*
623 *dependent*) protocol using a *producer-consumer* communication model and
624 providing a **unified architecture** throughout the manufacturing enterprise
625 [30][31].

626 CIP has been adapted in different types of network:

- 627 • **EtherNet/IP**, adaptation to *Transmission Control Protocol* (TCP) tech-
628 nologies
- 629 • **ControlNet**, adaptation to *Concurrent Time Domain Multiple Access*
630 (CTDMA) technologies
- 631 • **DeviceNet**, adaptation to *Controller Area Network* (CAN) technolo-
632 gies
- 633 • **CompoNet**, adaptation to *Time Division Multiple Access* (TDMA) tech-
634 nologies

635 **CIP objects** CIP is a *strictly object oriented* protocol at the upper layers:
636 each object of CIP has **attributes** (data), **services** (commands), **connec-
637 tions**, and **behaviors** (relationship between values and services of attributes)
638 which are defined in the **CIP object library**. The object library supports
639 many common automation devices and functions, such as analog and dig-
640 ital I/O, valves, motion systems, sensors, and actuators. So if the same
641 object is implemented in two or more devices, it will behave the same way
642 in each device [32].

643 **Security** [33] In EtherNet/IP implementation, security issues are the same
644 as in traditional Ethernet, such as network traffic sniffing and spoofing.
645 The use of the UDP protocol also exposes CIP to transmission route ma-
646 nipulation attacks using the *Internet Group Management Protocol* (IGMP)
647 and malicious traffic injection.

648 Regardless of the implementation used, it is recommended that certain
649 basic measures be implemented on the CIP network to ensure a high level
650 of security, such as *integrity, authentication and authorization*.

State of the Art

651 IN COVENTIONAL IT SYSTEMS, the objective of an attacker is to comprehend
652 the behavior of a program using diverse techniques in order to launch attacks that alter its execution flow, functionalities, or bypass limitations imposed by software licensing. These attack techniques involve an initial examination of the program, consisting of *static analysis* (i.e., analyzing the software without running it) and *dynamic analysis* (i.e., analyzing the program while it is running).

653 The outcome of these two investigative techniques is the *reverse engineering* of the software, which serves the purpose of identifying vulnerabilities or bugs and subsequently strategizing an attack.

654 In the context of OT systems, the notion of *reverse engineering* is not limited to its conventional definition, but also includes the concept of **process comprehension**. This term, introduced by Green et al. [9], refers to gaining a comprehensive understanding of the underlying physical process.

655 There is limited literature available concerning the gathering and analysis of information related to the comprehension and operation of an Industrial Control System (ICS). In Section 3.1, we will provide a brief overview of the existing literature on this topic, and in the subsequent sections, we will specifically focus on one of the presented papers.

670 3.1 Literature on Process Comprehension

671 **Keliris and Maniatikos** The first approach presented in this section is by
672 Keliris and Maniatakos [13]: they present a methodology for au-
673 tomating the reverse engineering of ICS binaries based on a *modular*
674 *framework* (called ICSREF) that can reverse binaries compiled with
675 CODESYS [34], one of the most popular and widely used PLC com-
676 pilers, irrespective of the language used.

677 **Yuan et al.** Yuan et al. [35] propose a *data-driven* approach to discover-
678 ing cyber-physical systems process behavior from data directly: to
679 achieve this goal, they have implemented a framework whose pur-
680 pose is to identify physical systems and transition logic inference,
681 and to seek to understand the mechanisms underlying these pro-
682 cesses, making furthermore predictions concerning their state trajec-
683 tories based on the discovered models.

684 **Feng et al.** Feng et al. [36] developed a framework that can generate sys-
685 tem *invariant rules* based on machine learning and data mining tech-
686 niques from ICS operational data log. These invariants are then se-
687 lected by systems engineers to derive IDS systems from them.

688 The experiment results on two different testbeds, the *Water Distri-*
689 *bution system* (WaDi) and the *Secure Water Treatment system* (SWaT),
690 both located at the iTrust - Center for Research in Cyber Security at
691 the University of Singapore of Technology and Design [37], show
692 that under the same false positive rate invariant-based IDSs have a
693 higher efficiency in detecting anomalies than IDS systems based on
694 a residual error-based model.

695 **Pal et al.** Pal et al. [38] work is somewhat related to Feng et al.'s: this
696 paper describes a data-driven approach to identifying invariants au-
697 tomatically using *association rules mining* [39] with the aim of generat-
698 ing invariants sometimes hidden from the design layout. The study
699 has the same objective of Feng et al.'s and uses too the iTrust SwaT

700 System as testbed.

701 Currently this technique is limited to only pair wise sensors and
702 actuators: for more accurate invariants generation, the technique
703 adopted must be capable of deriving valid constraints across multiple
704 sensors and actuators.

705 **Winnicki et al.** Winnicki et al. [11] instead propose a different approach
706 to process comprehension based on the *attacker's perspective* and not
707 limited to mere *Denial of Service* (DoS): their approach is to discover
708 the dynamic behavior of the system, in a semi-automated and process-
709 aware way, through *probing*, that is, slightly perturbing the cyber
710 physical system and observing how it reacts to changes and how
711 it returns to its original state. The difficulty and challenge for the
712 attacker is to perturb the system in such a way as to achieve an ob-
713 servable change, but at the same time avoid this change being seen
714 as a system anomaly by the IDSs.

715 **Green et al.** Green et al. [9] also adopt an approach based on the attacker's
716 perspective: this approach consists of two practical examples in a
717 *Man in the Middle* (MitM) scenario to obtain, correlate, and under-
718 stand all the types of information an attacker might need to plan an
719 attack to alter the process while avoiding detection.

720 The paper shows *step-by-step* how to perform a ICS **reconnaissance**, a
721 phase specifically designed to gather extensive intelligence on mul-
722 tiple fronts, including human factors, network and protocol infor-
723 mation, details about the manufacturing process, industrial applica-
724 tions, and potential vulnerabilities. The primary goal is to accumu-
725 late a wealth of information to enhance understanding and aware-
726 ness in these areas [40]).

727 Reconnaissance phase is fundamental to process comprehension and
728 thus to the execution of MitM attacks.

729 **Ceccato et al.** Ceccato et al. [12] propose a methodology based on a *black*
730 *box dynamic analysis* of an ICS using a reverse engineering tool to

32 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

731 derive from the scans performed on the memory registers of the ex-
732 posed PLCs and the Modbus protocol network scans an approximate
733 model of the physical process. This model is obtained by inferring
734 statistical properties, business process and system invariants from
735 data logs.

736 The proposed methodology was tested on a non-trivial case study,
737 using a virtualized testbed inspired by an industrial water treatment
738 plant.

739 In the next section we will examine this latest work in more detail,
740 which will be the basis for my work and thus the subsequent chap-
741 ters of this thesis.

742 3.2 Ceccato et al.'s black-box dynamic analysis 743 for water-tank systems

744 As previously mentioned, the paper introduces a methodology that re-
745 lies on black box dynamic analysis of an Industrial Control System (ICS)
746 and more particularly of its OT network. This methodology involves iden-
747 tifying potential Programmable Logic Controllers (PLCs) within the net-
748 work and scanning the memory registers of these identified controllers.
749 The purpose of this process is to obtain an approximate model of the con-
750 trolled physical process.

751 The primary goal of this black box analysis is to establish a correlation
752 between the different memory registers of the targeted PLCs and funda-
753 mental concepts of an OT network such as sensor values (i.e., measure-
754 ments), actuator commands, setpoints (i.e., range of values of a physical
755 variable), network communications, among others.

756 To accomplish this, the various types of memory registers are analyzed,
757 and attempts are made to determine the nature of the data they might
758 contain.

759 The second goal is to establish a relationship between the dynamic evo-
760 lution of these fundamental concepts.

761 To accomplish this, Ceccato et al. have developed a prototype tool [41]
762 that facilitates the reverse engineering of the physical system. This tool
763 goes through four distinct phases:

- 764 1. **scanning of the system and data pre-processing:** this phase involves
765 gathering data to generate data logs for the registers of PLCs and for
766 Modbus network communications.
- 767 2. **graphs and statistical analysis:** The collected data is utilized to pro-
768 vide insights into the memory registers associated with the Modbus
769 protocol by leveraging graphs and statistical data. This analysis ap-
770 proach offers valuable information about the characteristics and pat-
771 terns of the memory registers.
- 772 3. **invariants inference and analysis:** generates system invariants, which
773 are used to identify specific patterns and regularities within the sys-
774 tem. Additionally, this phase provides users with the capability to
775 view invariants related to a particular sensor or actuator.
- 776 4. **business process mining and analysis:** Using event logs, this phase
777 involves reconstructing the business process that depicts how a pro-
778 cess is executed. This step enables a thorough understanding of the
779 sequence of events that occur in the system and how they are in-
780 terrelated, ultimately leading to a comprehensive overview of the
781 business process.

782 Figure 3.1 presents a schematic representation of the stages and the
783 workflow associated with this work, specifying tools and technologies
784 used. In the subsequent sections of this chapter, we will provide a detailed
785 exploration of each of these phases, offering a comprehensive understand-
786 ing of the entire process.

34 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

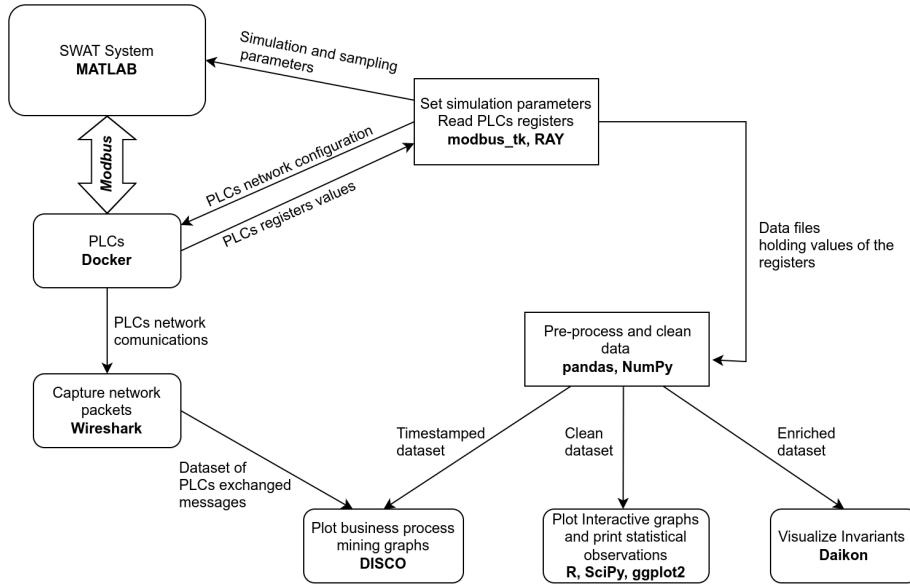


Figure 3.1: Workflow of Ceccato et al.'s stages and operations with used tools

787 3.2.1 Testbed

788 Before delving into the description of the methodology's different phases,
 789 let's first examine the testbed utilized to evaluate this approach. The testbed
 790 employed for testing purposes is a (very) simplified rendition of the iTrust
 791 SWaT system [15], as implemented by Lanotte et al. [42]. Figure 3.2 pro-
 792 vides a graphical representation of the testbed. This simplified version
 793 comprises three stages, each governed by a dedicated PLC.

794 **Stage 1** During the initial stage, a **tank** referred to as T-201 with a capacity
 795 of 80 gallons is filled with raw water using the P-101 pump. Connected to the T-201 tank, the MV-301 motorized valve flushes out the
 796 accumulated water from the tank, directing it to the next stage. Initially,
 797 the water flows from the T-201 tank to the *filtration unit* (which
 798 is not specifically identified by any sensor), and subsequently to a
 800 **second tank** denoted as T-202, with a capacity of 20 gallons.

801 **Stage 2** At the second stage, the water stored in tank T-202 flows into the
 802 *reverse osmosis unit* (RO), which serves as both a valve and a continu-

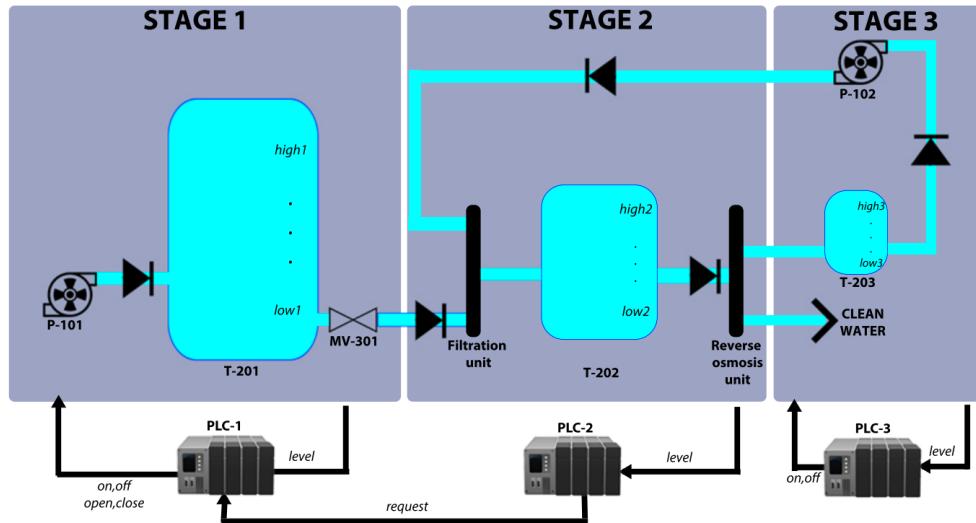


Figure 3.2: The simplified SWaT system used for running Ceccato et al. methodology

ous water extractor. The purpose of the RO unit is to reduce organic impurities present in the water. Subsequently, the water flows from the *RO unit* to the third and final stage of the system.

Stage 3 At the third stage, the water coming from the *RO unit* undergoes division based on whether it meets the required standards. If the water is deemed clean and meets the standards, it is directed into the distribution system. However, if the water fails to meet the standards, it is redirected to a *backwash tank* identified as T-203, which has a capacity of one gallon. The water stored in this tank is then pumped back to the stage 2 *filtration unit* using pump P-102.

As previously mentioned, each stage of the system is handled via a dedicated PLC, namely PLC1, PLC2, and PLC3, which are responsible for controlling their respective stages. Let's briefly explore the behavior of each PLC:

PLC1 PLC1 monitors the level of tank T-201 and distinguishes three different cases based on the level readings:

1. when the level of tank T-201 reaches the defined *low setpoint*

36 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

820 $low1$ (which is hardcoded in a specific memory register), PLC1
821 **opens pump P-101 and closes valve MV-301.** This configura-
822 tion allows the tank to be filled with water;

- 823 2. if the level of T-201 reaches the $high\ setpoint\ high1$ (which is also
824 hardcoded in a specific memory register), then the pump **P-101**
825 **is closed**;
- 826 3. in cases where the level of T-201 is between the $low\ setpoint\ low1$
827 and the $high\ setpoint\ high1$, PLC1 waits for a request from PLC2
828 to open or close the valve MV-301. If a request to open the valve
829 MV-301 is received, water will flow from T-201 to T-202. How-
830 ever, if no request is received, the valve remains closed. In both
831 situations, the pump P-101 remains closed.

832 **PLC2** PLC2 monitors the level of tank T-202 and adjusts its behavior based
833 on the water level. There are three cases to consider:

- 834 1. when the water level in tank T-202 reaches $low\ setpoint\ low2$ (also
835 hardcoded in the memory registers), PLC2 sends a request to
836 PLC1 through a Modbus channel to **open valve MV-301**. This
837 request is made in order to allow the water to flow from tank T-
838 201 to tank T-202. The transmission channel between the PLCs
839 is established by copying a boolean value from a memory reg-
840 ister of PLC2 to a corresponding register of PLC1.
- 841 2. when the water level in tank T-202 reaches the $high\ setpoint\ high2$
842 value (also hardcoded in the memory registers), PLC2 sends a
843 **close request to PLC1 for valve MV-301**. This request prompts
844 PLC1 to close the valve, stopping the flow of water from tank
845 T-201 to tank T-202.
- 846 3. In cases where the water level in tank T-202 is between the low
847 and high setpoints, the valve MV-301 remains in its current state
848 (open or closed) while the tank is either filling or emptying.

849 **PLC3** PLC3 monitors the level of the T-203 backwash tank and adjusts its
850 behavior accordingly. There are two cases to consider:

- 851 1. If the water level in the backwash tank reaches the *low setpoint*
852 *low3*, **PLC3 sets pump P103 to off**. This allows the backwash
853 tank to be filled.
- 854 2. If the water level in the backwash tank reaches the *high setpoint*
855 *high3*, **PLC3 opens pump P103**. This action triggers the pumping
856 of the entire content of the backwash tank back to the filter
857 unit of T-202.

858 **3.2.2 Scanning of the System and Data Pre-processing**

859 **Scanning tool** The Ceccato et al. scanning tool extends and generalizes
860 a project I did [43] for the "Network Security" and "Cyber Security for IoT"
861 courses taught by Professors Massimo Merro and Mariano Ceccato, re-
862 spectively, in the 2020/21 academic year. The original project involved,
863 in its first part, the recognition within a network of potential PLCs lis-
864 tening on the standard Modbus TCP port 502 using the Nmap module
865 for Python, obtaining the corresponding IP addresses: then a (sequential)
866 scan of a given range of the memory registers of the found PLCs was per-
867 formed to collect the register data. The data thus collected were saved to
868 a file in *JavaScript Object Notation* (JSON) format for later use in the second
869 part of my project.

870 The scanning tool by Ceccato et. al works in a similar way, but extends
871 what originally did by trying to discover other ports on which the Mod-
872 bus protocol might be listening (since in many realities Modbus runs on
873 different ports than the standard one, according to the concept of *security*
874 *by obscurity*) and, most importantly, by **parallelizing and distributing the**
875 **scan** of PLC memory registers through the Ray module [44], specifying
876 moreover the desired granularity of the capture. An example of raw data
877 capture can be seen at Listing 3.1:

```
878        "127.0.0.1/8502/2022-05-03 12_10_00.591": {  
879            "DiscreteInputRegisters": {"%IX0.0": "0"},  
880            "InputRegisters": {"%IW0": "53"},  
881            "HoldingOutputRegisters": {"%QW0": "0"},
```

38 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

```
882     "MemoryRegisters": {"%MW0": "40", "%MW1": "80"},  
883     "Coils": {"%QX0.0": "0"}}
```

Listing 3.1: Example of registers capture

884 The captured data includes PLC's IP address, Modbus port and timestamp
885 (first line), type and name of registers with their values read from the scan
886 (subsequent lines).

887 The tool furthermore offers the possibility, in parallel to the memory
888 registers scan, of **sniffing network traffic** related to the Modbus protocol
889 using the *Man in the Middle* (MitM) technique on the supervisory control
890 network using a Python wrapper for tshark/Wireshark [45] [46]. An ex-
891 ample of raw data obtained with this sniffing can be seen in Listing 3.2:

```
892     Time, Source, Destination, Protocol, Length, Function Code,  
893     ↪ Destination Port, Source Port, Data, Frame length on the  
894     ↪ wire, Bit Value, Request Frame, Reference Number, Info  
895     2022-05-03 11:43:58.158, IP_PLC1, IP_PLC2, Modbus/TCP, 76, Read  
896     ↪ Coils, 46106, 502, , 76, TRUE, 25, , "Response: Trans: 62;  
897     ↪ Unit: 1, Func: 1: Read Coils"
```

Listing 3.2: Example of raw network capture

898 **Data Pre-processing** The data collected by scanning the memory regis-
899 ters of the PLCs are then reprocessed by a Python script and converted
900 in order to create a distinct raw dataset in *Comma Separated Value* for-
901 mat (CSV) for each PLC, containing the memory register values associ-
902 ated with the corresponding controller registers. These datasets are repro-
903 cessed again through the Python modules for **pandas** [47] and **NumPy** [48]
904 by another script to first perform a **data cleanup**, removing all unused reg-
905 isters, **merged** into a single dataset, and finally **enriched** with additional
906 data, such as the **previous value** of all registers and the **measurement**
907 **slope**, that is, the trend of the water level in the system tanks along the
908 system cycles¹. See 3.2.7 for more detail.

¹Not all additional data are calculated and entered automatically by the tool: some are manually inserted.

909 This process leads to the creation of two copies of the full dataset: one
 910 enriched with the additional data, but not timestamped, which will be
 911 used for the invariant analysis; the other unenriched, but timestamped,
 912 which will be used for business process mining.

913 **3.2.3 Graphs and Statistical Analysis**

914 The paper mentions the presence of a *mild graph analysis*, performed
 915 using the framework **R** [49] for statistical analysis at the time of data gath-
 916 ering to find any uncovered patterns, trends and identify measurements
 917 and/or actuator commands through the analysis of registers holding mu-
 918 table values.

919 There is actually no trace of this within the tool: *graph analysis* and *sta-*
 920 *tistical analysis* of the data contained in the PLC memory registers are in-
 921 stead performed using the **matplotlib libraries** and statistical algorithms
 922 made available by the **SciPy libraries** [50], through two separate Python
 923 scripts (see Figure 3.3).

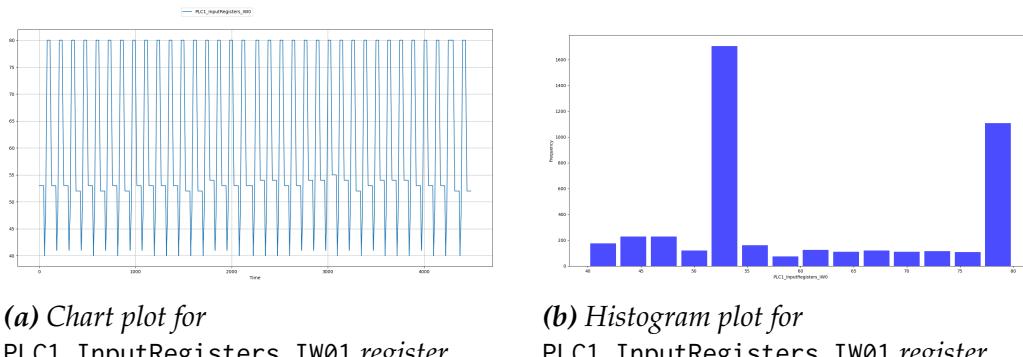


Figure 3.3: Output graphs from graph analysis

924 The first script plots the charts, one at the time, of certain registers en-
 925 tered by the user from the command line, plots in which one can see the
 926 trend of the data and get a first basic idea of what that particular regis-
 927 ter contains (a measurement, an actuation, a hardcoded setpoint, ...) and
 928 possibly the trend; the second script, instead, shows a **histogram and sta-**

40 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

929 **tistical informations** about the register entered as command-line input.
930 These informations include:

- 931 • the mean, median, standard deviation, maximum value and mini-
932 mum value
- 933 • two tests for the statistical distribution: *Chi-squared* test for unifor-
934 mity and *Shapiro-Wilk* test for normality, as shown in Listing 3.3:

```
935     Chi-squared test for uniformity
936     Distance      pvalue      Uniform?
937     12488.340    0.00000000    NO
938
939     Shapiro-Wilk test for normality
940     Test statistic   pvalue      Normal?
941     0.844      0.00000000    NO
942
943     Stats of PLC1_InputRegisters_IW0
944     Sample mean = 60.8881; Stddev = 13.0164; max = 80; min =
945     ↪ 40 for 4488 values
```

Listing 3.3: Statistical data for PLC1_InputRegisters_IW0 register

946 3.2.4 Invariant Inference and Analysis

947 For invariant analysis Ceccato et al. rely on **Daikon** [14], a framework
948 to **dynamically detect likely invariants** within a program. An *invariant*
949 is a property that holds at one or more points in a program, properties
950 that are not normally made explicit in the code, but within assert state-
951 ments, documentation and formal specifications: invariants are useful in
952 understanding the behavior of a program (in our case, of the cyber physi-
953 cal system).

954 Daikon uses *machine learning* techniques applied to arbitrary data with
955 the possibility of setting custom conditions for analysis by using a spe-
956 cific file [51] with a *.spinfo* extension (see Listing 3.4). The framework is
957 designed to find the invariants of a program, with various supported pro-
958 gramming languages, starting from the direct execution of the program

959 itself or passing as input the execution run (typically a file in CSV format);
 960 the authors of the paper tried to apply it by analogy also to the execution
 961 runs of a cyber physical system, to extract the invariants of this system.

```
962 PPT_NAME aprogram.point:::POINT
963 VAR1 > VAR2
964 VAR1 == VAR3 && VAR1 != VAR4
```

Listing 3.4: Generic example of a .spininfo file for customizing rules in Daikon

965 Therefore, Daikon is fed with the enriched dataset obtained in the pre-
 966 processing phase²: a simple bash script launches Daikon (optionally spec-
 967 ifying the desired condition for analysis in the *.spininfo* file), which output is
 968 simply redirected to a text file containing the general invariants of the sys-
 969 tem (i.e., valid regardless of any custom condition specified), those gener-
 970 ated based on the custom condition in the *.spininfo* file, and those generated
 971 based on the negation of the condition (see Listing 3.5 below).

```
972 =====
973 aprogram.point:::POINT
974 PLC2_MemoryRegisters_MW1 == PLC3_MemoryRegisters_MW1
975 PLC1_MemoryRegisters_MW0 == 40.0
976 PLC1_MemoryRegisters_MW1 == 80.0
977 PLC1_Coils_QX00 one of { 0.0, 1.0 }
978 [...]
979 =====
980 aprogram.point:::POINT; condition="PLC1_InputRegisters_IW0
981   ↪ > 60"
982 PLC1_InputRegisters_IW0 > PLC1_MemoryRegisters_MW0
983 PLC1_InputRegisters_IW0 > PLC1_Min_safety
984 PLC1_MemoryRegisters_MW0 < prev_PLC1_InputRegisters_IW0
985 [...]
986 =====
987 aprogram.point:::POINT; condition="not(
988   ↪ PLC1_InputRegisters_IW0 > 60)"
989 PLC1_InputRegisters_IW0 < PLC1_MemoryRegisters_MW1
990 PLC1_InputRegisters_IW0 < PLC1_Max_safety
```

²In the paper, timestamped dataset is explicitly mentioned as input: from the tests performed, Daikon seems to ignore timestamps, hence it is indifferent whether the dataset is timestamped or not

42 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

```
991     PLC1_MemoryRegisters_MW1 > prev_PLC1_InputRegisters_IW0  
992     [...]
```

Listing 3.5: The three sections of Daikon analysis outcomes

993 When the analysis is finished, the user is asked to enter the name of a reg-
994 istry to view its related invariants.

995

996 Some examples of invariants derived from the enriched dataset may be:

- 997 • measurements bounded by some setpoint;
- 998 • actuators state changes occurred in the proximity of setpoints or,
999 vice versa, proximity of setpoints upon the occurrence of an actuator
1000 state change;
- 1001 • state invariants of some actuators correspond to a specific trend in
1002 the evolution of the measurements (ascending, descending, or sta-
1003 ble) or, vice versa, the measurements trend corresponds to a specific
1004 state invariant of some actuators.

1005 3.2.5 Business Process Mining and Analysis

1006 *Process mining* is the analysis of operational processes based on the
1007 event log [52]: the aim of this analysis is to **extract useful informations**
1008 from the event data to **reconstruct and understand the behavior** of the
1009 business process and how it was actually performed.

1010 In the considered system, process mining begins by analyzing the event
1011 logs derived from scanning the memory registers of the PLCs and moni-
1012 toring the network communications associated with the Modbus protocol,
1013 as detailed in Subsection 3.2.2. These event logs serve as the *execution trace*
1014 of the system. A Java program is utilized to extract and consolidate infor-
1015 mation from these event logs, resulting in a CSV format file that captures
1016 the relevant data.

1017 This file is fed to **Disco** [53], a commercial process mining tool, which
1018 generates an *activity diagram* similar to UML Activity Diagram and whose

nodes represent the activities while the edges represent the relations between these activities. In Figure 3.4 we can see an example of this diagram referred to PLC2 of the testbed: nodes represent the trend of register associated with measurement, actuator state changes, and communications between PLCs involving these state changes, while edges represent transitions with their associated time duration and frequency.

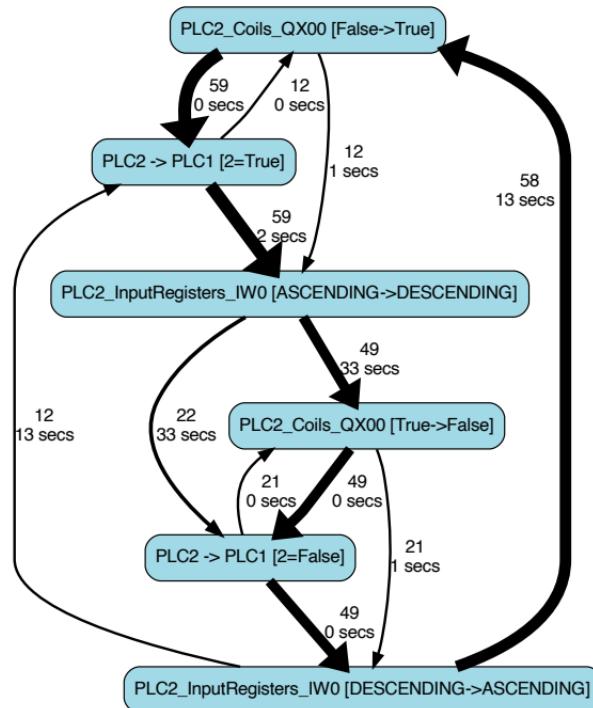


Figure 3.4: An example of Disco generated activity diagram for PLC2

The *business process* obtained in this way provides an **overview of the system** and makes it possible to **make conjectures** about its behavior, particularly between changes in actuator state and measurement trends (i.e., a given change in state of some actuators corresponds to a specific measurement trend and vice versa), and with the possibility of **establishing causality** between Modbus communications and state changes within the physical system.

44 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

1032 3.2.6 Application

1033 In this section we will see how the black box analysis presented above
1034 in its various phases is applied in practice, using the testbed described in
1035 Subsection 3.2.1. The methodology supports a *top-down approach*: that
1036 is, we start with an overview of the industrial process and then gradually
1037 refine our understanding of the process by descending to a higher and
1038 higher level of detail based on the results of the previous analyses and
1039 focusing on the most interesting parts of the system for further in-depth
1040 analysis.

1041 **Data Collection and Pre-processing** According to what is described in
1042 the paper, the data gathering process lasted six hours, with a granular-
1043 ity of one data point per second (a full system cycle takes approximately
1044 30 minutes). Each datapoint consists of 168 attributes (55 registers plus
1045 a special register concerning the tank slope of each PLC) after the en-
1046 richment. In addition, IP addresses are automatically replaced by an ab-
1047 stract name identified by the prefix PLC followed by a progressive integer
1048 (PLC1, PLC2, PLC3), in order to make reading easier.

1049 **Graphs and Statistical Analysis** Graphs and Statistical Analysis revealed
1050 three properties regarding the contents of the registers:

1051 **Property 1:** PLC1_MemoryRegisters_MW0, PLC1_MemoryRegisters_MW1,
1052 PLC2_MemoryRegisters_MW0, PLC2_MemoryRegisters_MW1,
1053 PLC3_MemoryRegisters_MW0 and PLC3_MemoryRegisters_MW1
1054 registers contain constant integer values (40, 80, 10, 20, 0, 10 respec-
1055 tively)³. The authors speculate that they may be (relative) hardcoded
1056 **setpoints**.

³From my tests on the original tool and dataset, the PLC3_MemoryRegisters_MW0 register is deleted during the *pre-processing* phase, as it is recognized as an unused register because of the constant value "0" it takes on. This leads me to assume that the properties are derived from a human read of the dataset prior to the *pre-processing* phase.

1057 Property 2: PLC1_Coils_QX01, PLC1_Coils_QX02, PLC2_Coils_QX01,
 1058 PLC2_Coils_QX02, PLC3_Coils_QX01 and PLC3_Coils_QX03 contain mu-
 1059 table binary (Boolean) values. The authors speculate that these reg-
 1060 isters can be associated with the **actuators** of the system.

1061 Property 3: PLC1_InputRegisters_IW0, PLC2_InputRegisters_IW0 and
 1062 PLC3_InputRegisters_IW0 registers contain mutable values.

1063 Property 3 suggests that those registers might contain **values related to**
 1064 **measurements**: it is therefore necessary to investigate further to see if the
 1065 conjecture (referred to as *Conjecture 1* in the paper) is correct.

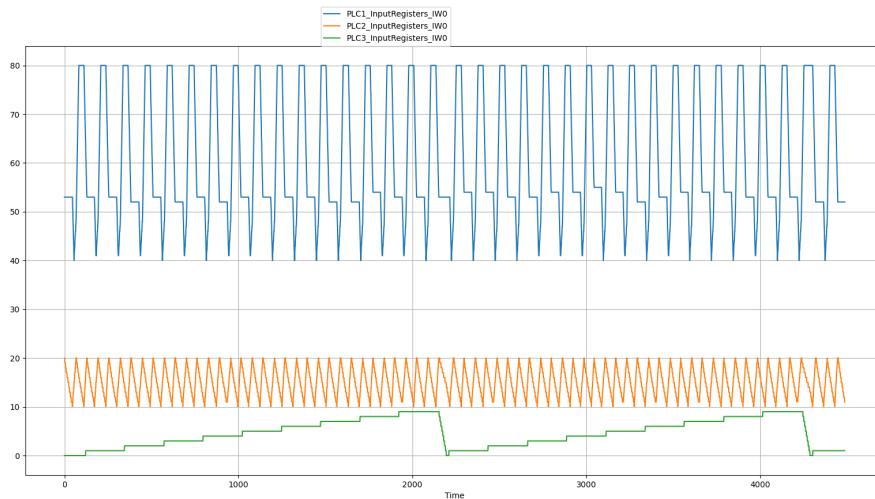


Figure 3.5: Execution traces of *InputRegisters_IW0* on the three PLCs

1066 The graph analysis of the *InputRegisters_IW0* registers of the three
 1067 PLCs (summarized in Figure 3.5 with a single plot) not only seems to con-
 1068 firm the conjecture, but also allows the measurements to be correlated with
 1069 the contents of the *MemoryRegisters_MW0* and *MemoryRegisters_MW1* regis-
 1070 ters to the measurements, which may well represent the **relative setpoints**
 1071 **of the measurements**. Hence, we have *Conjecture 2* described in the paper
 1072 referring to the relative setpoints:

1073

1074 *Conjecture 2:*

1075 - the relative setpoints for *PLC1_InputRegisters_IW0* are 40 and 80;

46 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

- 1076 - the relative setpoints for PLC2_InputRegisters_IW0 are 10 and 20;
1077 - the relative setpoints for PLC3_InputRegisters_IW0 0 and 9.

1078 Further confirmation of this conjecture may come from statistical anal-
1079 ysis. Indeed, in the example in Listing 3.1, some statistical data are given
1080 for the register PLC1_InputRegisters_IW0, including the maximum value
1081 and the minimum value: these values are, in fact, 80 and 40 respectively.

1082 **Business Process Mining and Analysis** With Business Process Mining,
1083 the authors aim to **visualize and highlight relevant system behaviors** by
1084 relating PLC states and Modbus commands.

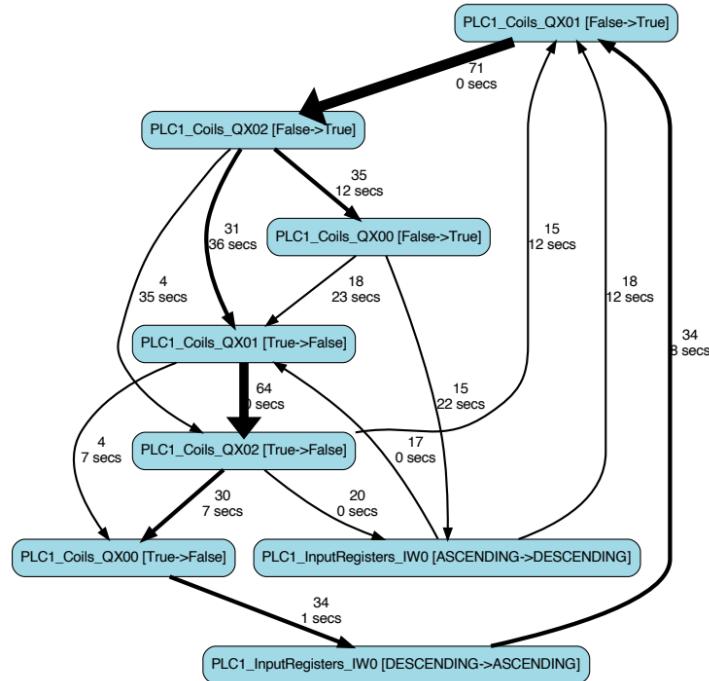
1085 Through analysis of the activity diagrams shown in Figure 3.6, drawn
1086 through Disco, they derive the following properties and conjectures:

1087 **Property 4:** PLC2 sends messages to PLC1 (see Figure 3.6b) which are
1088 recorded to PLC1_Coils_QX02.

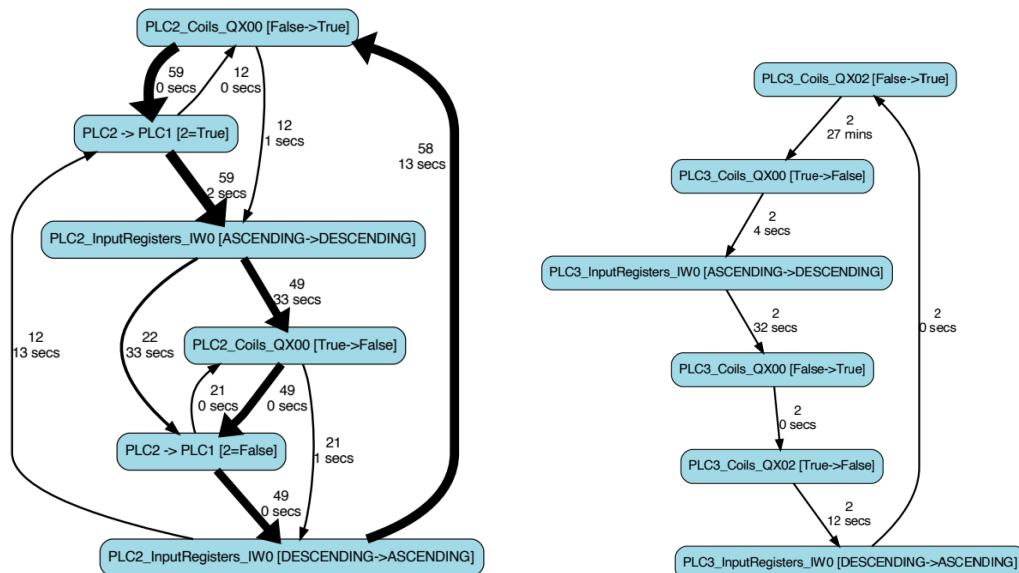
1089 **Conjecture 3:** PLC2_Coils_QX00 determines the trend in tank T-202 (Figure
1090 3.6b).
1091 When this register is set to *True*, the input register PLC2_InputRegisters_IW0
1092 related to the tank controlled by PLC2 starts an **ascending trend**; vice
1093 versa, when the coil register is set to *False*, the input register starts a
1094 **descending trend**.

1095 **Conjecture 4:** If PLC1_Coils_QX00 change his value to True, trend in tank
1096 T-201, related to PLC1_InputRegisters_IW0 and controlled by PLC1,
1097 become **ascending** (see Figure 3.6a)

1098 **Conjecture 5:** PLC3_Coils_QX00 starts a **decreasing trend** in tank T-203, re-
1099 lated to PLC3_InputRegisters_IW0 and controlled by PLC3, whereas
1100 PLC3_Coils_QX02 starts an **increasing trend** on the tank (see Figure
1101 3.6c)



(a) States in PLC1



(b) States and Modbus command in PLC2

(c) States in PLC3

Figure 3.6: Business process with states and Modbus commands for the three PLCs

48 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

1102 **Invariant Inference and Analysis** The last phase of the analysis of the
1103 example industrial system is invariant analysis, performed through Daikon
1104 framework. At this stage, an attempt will be made to confirm what has
1105 been seen previously and to derive new properties of the system based on
1106 the results of the Daikon analysis.

1107 To get gradually more and more accurate results, the authors presum-
1108 ably performed more than one analysis with Daikon, including certain
1109 rules within the *splitter information file* (see Section 3.2.4 and Listing 3.4)
1110 based on specific conditions placed on the measurements, for example, the
1111 level of water contained in a tank. Given moreover the massive amount
1112 of invariants generated by Daikon's output, it is not easy to identify and
1113 correlate those that are actually useful for analysis: this must be done man-
1114 ually.

1115 However, it was possible to have confirmation of the conjectures made
1116 in the previous stages of the analysis: starting with the setpoints, analyz-
1117 ing the output of the invariants returned by Daikon⁴ reveals that

1118
1119 PLC1_InputRegisters_IW0 >= PLC1_MemoryRegisters_MW0 == 40.0
1120 PLC1_InputRegisters_IW0 <= PLC1_MemoryRegisters_MW1 == 80.0
1121 PLC2_InputRegisters_IW0 >= PLC2_MemoryRegisters_MW0 == 10.0
1122 PLC2_InputRegisters_IW0 <= PLC2_MemoryRegisters_MW1 == 20.0
1123 PLC3_InputRegisters_IW0 >= PLC3_MemoryRegisters_MW0 == 0.0
1124 PLC3_InputRegisters_IW0 <= PLC3_MemoryRegisters_MW1 == 9.0
1125
1126 i.e., that the MemoryRegisters_MW0 and MemoryRegisters_MW1 registers of
1127 each PLC contain the **absolute minimum and maximum setpoints**, re-
1128 spectively (*Property 5*).

1129 There is also a confirmation regarding *Property 4*: from the computed
1130 invariants it can be seen that

⁴The invariants shown here are a manual summary and derivation of those actually returned in output by Daikon. We will discuss this more in Section 3.2.7

1131
1132 PLC1_Coils_QX01 == PLC1_Coils_QX02 == PLC2_Coils_QX00
1133

1134 and from this derive that there is a **communication channel between PLC2**
1135 **and PLC1**, where the value of PLC2_Coils_QX00 is copied to PLC1_Coils_QX01
1136 and PLC1_Coils_QX02 (*Property 6*).

1137 Regarding the **relationships between actuator state changes and mea-**
1138 **surement trends**, invariant analysis yields the results summarized in the
1139 following rules:

1140 **Property 7:** Tank T-202 level *increases* iff PLC1_Coils_QX01 == True. Oth-
1141 erwise, if PLC1_Coils_QX01 == False will be *non-increasing*.

1142 This is because if the coil is *True* the condition

1143 PLC2_InputRegisters_IW0 == PLC2_MemoryRegisters_MW0 == 20.0 && PLC2_slope > 0
1144 is verified. On the opposite hand, if the coil is *False*, the condition
1145 PLC2_InputRegisters_IW0 == PLC2_MemoryRegisters_MW0 == 20.0 && PLC2_slope <= 0 is verified. The
1146 *slope* is increasing if > 0, decreasing if < 0, stable otherwise.

1147 **Property 8:** Tank T-201 level *increases* iff PLC1_Coils_QX00 == True. On the
1148 other hand, if PLC1_Coils_QX00 == False and if PLC1_Coils_QX01 ==
1149 True the level will be *non-decreasing*.

1150 **Property 9:** Tank T-203 level *decreases* iff PLC3_Coils_QX00 == True. It will
1151 be *non-decreasing* if PLC1_Coils_QX00 == False.

1152 The last two properties concern the **relationship between actuator state**
1153 **changes and the setpoints**: it is intended to check what happens to the
1154 actuators when the water level reaches one of these setpoints. From the
1155 analysis of the relevant invariants, the following properties are derived:

1156 **Property 10:** Tank T-201 reaches the upper absolute setpoint when
1157 PLC1_Coils_QX00 changes its state from *True* to *False*. If the coil changes
1158 from *False* to *True*, the tank reaches its absolute lower setpoint.

50 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

1159 **Property 11:** Tank T-203 reaches the upper absolute setpoint when
1160 PLC3_Coils_QX00 changes its state from *True* to *False*. If the coil changes
1161 from *False* to *True*, the tank reaches its absolute lower setpoint.

1162 3.2.7 Limitations

1163 The methodology proposed by Ceccato et al. is certainly valid and
1164 offers a good starting point for approaching the reverse engineering of
1165 an industrial control system from the attacker's perspective, while also
1166 providing a tool to perform this task.

1167 The limitations of this approach, however, all lie in the tool mentioned
1168 above and also in the testbed described in Section 3.2.1. In this section
1169 we will explain which are the criticisms of each phase, while in Chapter 4
1170 we will formulate proposals to improve and make this methodology more
1171 efficient.

1172 **General Criticism** There are several critical aspects associated with the
1173 application of this approach: the primary one concerns the fact that the
1174 proposed tool seems to be built specifically for the testbed used and that
1175 it is not applicable to other contexts, even to the same type of industrial
1176 control system (water treatment systems, in this case).

1177 What severely limits the analysis performed with the tool implemented
1178 by Ceccato et al. is the use of *ad hoc* solutions and *a posteriori* interventions
1179 done manually on the datasets after the data gathering process: we will
1180 discuss this last aspect in more detail later.

1181 Moreover, there is the presence of many *hardcoded* variables and condi-
1182 tions within the scripts: this makes the system unconfigurable and unable
1183 to properly perform the various stages of the analysis as errors can occur
1184 due to incorrect data and mismatches with the system under analysis.
1185 Having considered, furthermore, only the Modbus protocol for network
1186 communications between the PLCs is another major limiting factor and

1187 does not help the methodology to be adaptable to different systems com-
1188 municating with different protocols (sometimes even multiple ones on the
1189 same system).

1190 Let us now look at the limitations and critical aspects of each phase.

1191 **Testbed** The testbed environment used by Ceccato et. al is entirely simu-
1192 lated, from the physical system to the control system. The PLCs were built
1193 with **OpenPLC** [54] in a Docker environment [55], while the physics part
1194 was built through **Simulink** [56].

1195 OpenPLC is an open source cross-platform software that simulates the
1196 hardware and software functionality of a physical PLC and also offers a
1197 complete editor for PLC program development with support for all stan-
1198 dard languages: *Ladder Logic* (LD), *Function Block Diagram* (FBD), *Instruc-*
1199 *tion List* (IL), *Structured Text* (ST), and *Sequential Function Chart* (SFC).

1200 It is for sure an excellent choice for creating a zero-cost industrial or
1201 home automation and *Internet of Things* (IoT) system that is easy to man-
1202 age via a dedicated, comprehensive and functional web interface. In spite
1203 of these undoubted merits, however, there are (at the moment) **very few**
1204 **supported protocols**: the main one and also referred to in the official doc-
1205 umentation is **Modbus**, while the other protocol is DNP3.

1206 **First limitation** The biggest problem with the testbed, however, is not
1207 with the controller part, but with the **physical part**: first of all, it
1208 must be said that although this is something purely demonstrative
1209 even though it is fully functional, the implemented Simulink model
1210 is really **oversimplified** compared to other testbeds. In fact, in the
1211 entire system there are only three actuators, two of which are con-
1212 nected to the same tank and controlled by the same PLC, and sensors
1213 related only to the water level in the system's tanks: in a real sys-
1214 tem there are many more *field devices*, which can monitor and control
1215 other aspects of the system beyond the mere contents of the tanks.
1216 Consider, for example, measuring and controlling the chemicals in

52 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

1217 the water, the pressure of the liquid in the filter unit, or more simply
1218 the amount of water flow at a given point or time.

1219 All these must be considered and represent a number of additional
1220 variables that makes analysis and consequently reverse engineering
1221 of the system more difficult.

1222 ***Second limitation*** The second critical aspect concerns the **simulation of**
1223 **the physics of the liquid** inside the tanks: Simulink does not con-
1224 sider the fact that inside a tank that is filling (emptying) the liquid
1225 in it undergoes **fluctuations** which cause the level sensor not to see
1226 the water level constantly increasing (decreasing) or at most being
1227 stable at each point of detection. Figure 3.7 exemplifies more clearly
1228 with an example the concept just expressed: these oscillations cause
1229 a **perturbation** in the data.

1230 This issue leads to the difficulty, on a real physical system, of **cor-**
1231 **rectly calculating the trend of a measurement** by using the slope
1232 attribute: if this was obtained with a too low granularity, the trend
1233 will be oscillating between increasing and decreasing even when in
1234 reality this would be in general increasing (decreasing) or stable; on
1235 the other hand, if the slope was obtained with a too high granularity
1236 there is a loss of information and the trend may be "flattened" with
1237 respect to reality.

1238 In the present case, the slope in the Simulink model was calculated
1239 statically with a (very) low granularity, 5 and 6 seconds according
1240 to the Properties 7 and 9 described in the original paper: an aver-
1241 agely careful reader will have already guessed that this granularity
1242 is inapplicable to the real system in Figure 3.7b. As we will later see,
1243 we need to **operate on the data perturbations** to be able to obtain a
1244 suitable granularity and a correct calculation of the slope and conse-
1245 quently of the measurement trend.

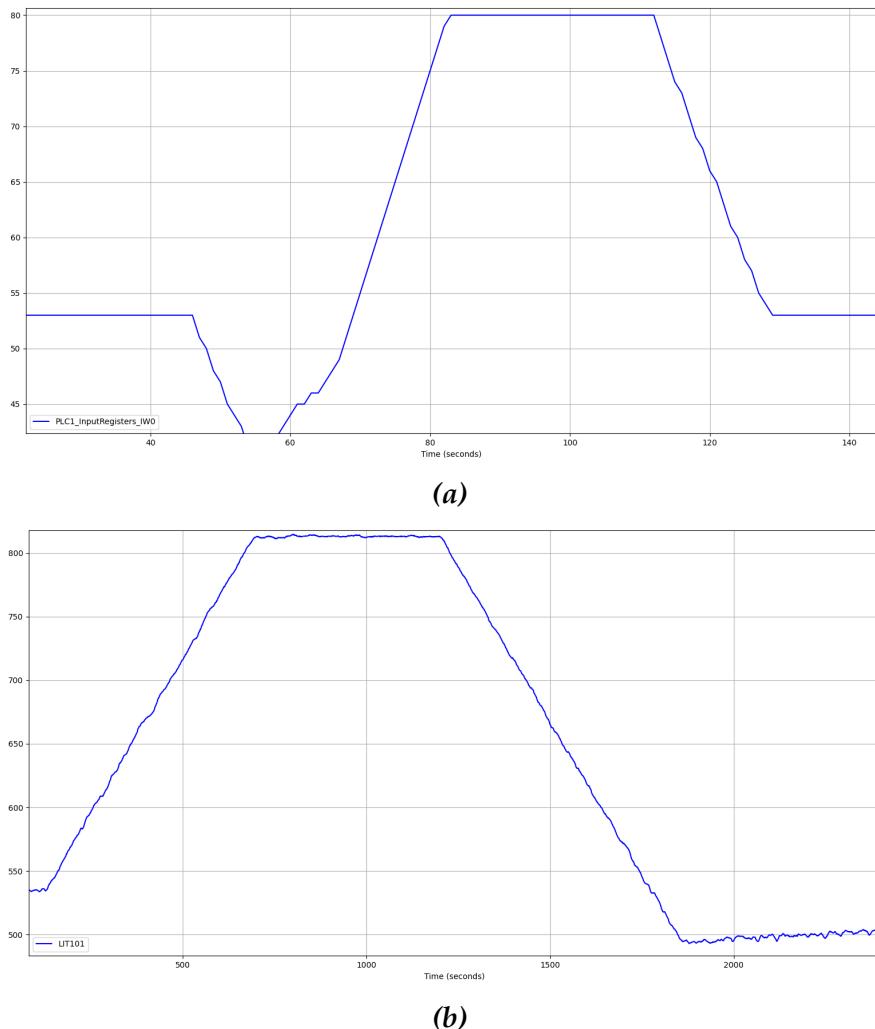


Figure 3.7: Water physics compared: simulated physics in the Simulink model (a) and physics in a real system (iTrust SWaT) (b). Fluctuations in the tank level in (b), almost completely absent in (a), can be appreciated.

1246 **Pre-processing** In the pre-processing phase, the authors make use of a
 1247 Python script to merge all the datasets of the individual PLCs into a single
 1248 dataset, remove the (supposedly) unused registers, and finally enrich the
 1249 obtained dataset with additional attributes. These attributes, as seen in
 1250 3.2.2, are:

- 1251 • the **previous value** of all registers;

54 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

- 1252 • some **additional relative setpoints** named PLC x _Max_safety and
1253 PLC x _Min_safety (where x is the PLC number), which represent a
1254 kind of alert on reaching the maximum and minimum water levels
1255 of the tanks;
- 1256 • the **measurement slope**.

1257 ***First limitation*** Merging the datasets of all individual PLCs into a single
1258 dataset representing the entire system can be a sound practice if the
1259 system to be analyzed is (very) small as is the testbed analyzed here,
1260 consisting of a few PLCs and especially a few registers. If, however,
1261 the complexity of the system increases, this type of merging can be-
1262 come counterproductive and make it difficult to analyze and under-
1263 stand the data obtained in subsequent steps.

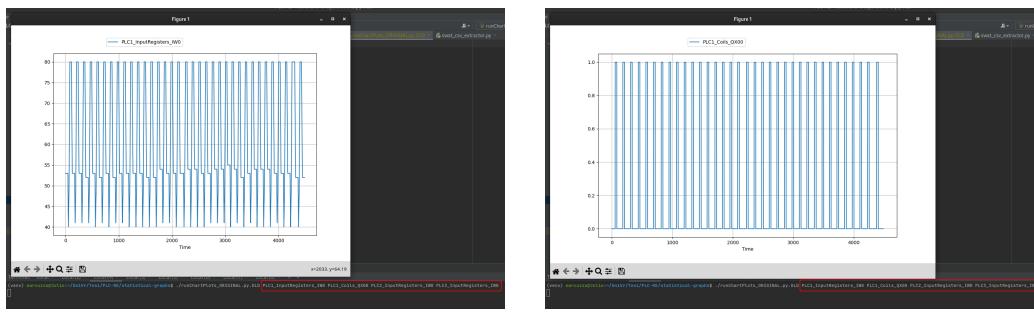
1264 In short, there is no possibility to analyze only a subsystem and thus
1265 make the analysis faster and more understandable. Moreover, a data
1266 gathering can take up to days, and the analyst/attacker may need to
1267 make an analysis of the system isolating precise time ranges, ignor-
1268 ing everything that happens before and/or after: all of this, with the
1269 tool we have seen, cannot be done.

1270 ***Second limitation*** Regarding the additional attributes, looking at the code
1271 of the script that performs the enrichment, we observed that **some at-**
1272 **tributes were manually inserted** after the merging phase: we are re-
1273 ferring in particular to the attributes PLC x _Max_safety and PLC x _Min_safety,
1274 whose references were moreover hardcoded into the script, and the
1275 *slope* whose calculation method we mentioned in the previous para-
1276 graph about the testbed limitations.

1277 In the end, only the attribute *prev* related to the value at the previous
1278 point of the detection is inserted automatically for all registers, more-
1279 over without the possibility to choose whether this attribute should
1280 be extended to all registers or only to a part.

1281 **Graphs and Statistical Analysis** Describing the behavior of graphical
1282 analysis in Section 3.2.3 we had already mentioned that only one register

1283 plot at a time was shown and not, for example, a single window containing
 1284 the charts of all registers entered by the user as input from the com-
 1285 mand line, such as in Figure 3.5. Figure 3.8 shows the actual behavior
 1286 of graphical analysis: note that although we have specified four registers
 1287 (highlighted in red in the figures) as command-line parameters, only one
 1288 at a time is shown and it is necessary to close the current chart in order to
 1289 display the next one.



(a) Chart for PLC1_InputRegisters_IW0

(b) Chart for PLC1_Coils_QX00

Figure 3.8: Behavior of the Graph Analysis on the Ceccato et al.'s tool

1290 **First limitation** While displaying charts for individual registers still pro-
 1291 vides useful information about the system such as the distinction
 1292 between actuators and measurements and the general trend of the
 1293 latter, single display does not allow one to catch, or at least makes it
 1294 difficult, the relationship that exists between actuators and measure-
 1295 ments, where it exists, because a view of the system as a whole is
 1296 missing.

1297 In this way, the risk is to make conjectures about the behavior of the
 1298 system that may prove to be at least imprecise, if not inaccurate.

1299 **Second limitation** On the other hand, regarding the statistical analysis,
 1300 two observations need to be made: the first is that for the given sys-
 1301 tem, I personally was unable to appreciate the usefulness of the gen-
 1302 erated histogram in Figure 3.3b, as it does not provide any particular
 1303 new information that has not already been obtained from the graph-

56 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

1304 ical analysis (except maybe something marginal); the second obser-
1305 vation pertains to the presentation of statistical information obtained
1306 from the histogram plot. In certain cases, the histogram plot itself
1307 can overshadow the displayed statistical information. These statis-
1308 tics are actually shown on the terminal from which the script is exe-
1309 cuted. However, to an inattentive or unfamiliar user, these statistics
1310 may be mistaken for debugging output or warnings, as they coin-
1311 cide with the display of the histogram plot window, which takes the
1312 focus (see Figure 3.9).

1313

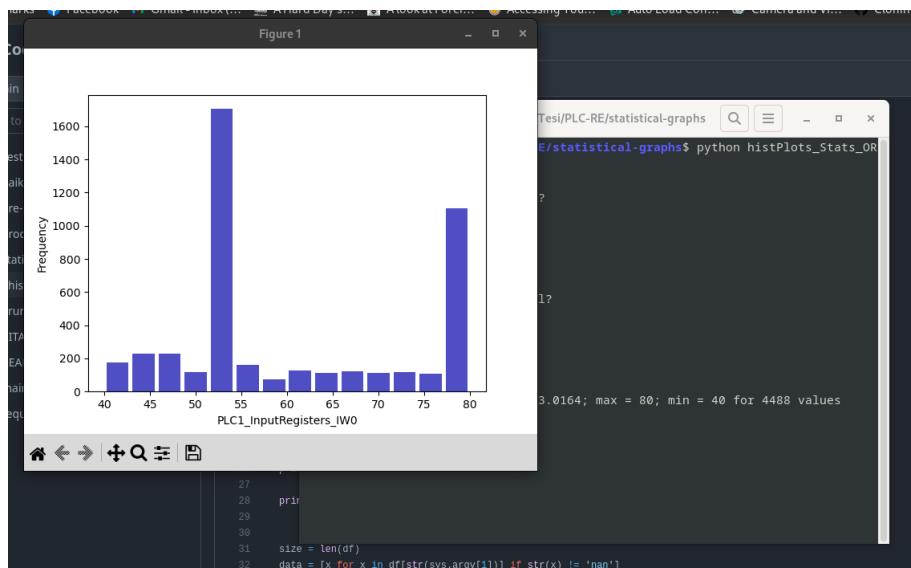


Figure 3.9: Histogram plot overshadowing statistical information shown on the terminal window in the background

1314

In general, however, little statistical information is provided.

1315 **Business Process Mining and Analysis** Concerning the data mining,
1316 this is a purely *ad hoc solution*, designed to work under special conditions:
1317 first, the timestamped dataset of the physical process and the one obtained
1318 after the packet sniffing operation of Modbus traffic on the network need
1319 to be synchronized and have the same granularity, in this case one event
1320 per second.

1321 It is relatively easy, therefore, to find correspondences between Mod-
1322 bus commands sent over the network and events occurring on the physical
1323 system, such as state changes in actuations, due in part to the fact that the
1324 number of communications over the network is really small (see Section
1325 3.2.1).

1326 ***First limitation*** In a real system, network communications are much more
1327 numerous and involve many more devices even in the same second:
1328 finding the exact correspondence with what is happening in the cy-
1329 ber physical system becomes much more difficult.

1330 Since this is, as mentioned, an *ad hoc* solution, only the Modbus pro-
1331 tocol is being considered: as widely used as this industrial protocol
1332 is, other protocols that are widely used [57] such as EtherNet/IP (see
1333 Section 2.2.6.2) or Profinet should be considered in order to extend
1334 the analysis to other industrial systems that use a different commu-
1335 nication network.

1336 ***Second limitation*** The other limiting aspect of the business process min-
1337 ing phase is the **process mining software** used to generate the ac-
1338 tivity diagram. As mentioned in Section 3.2.5, the process mining
1339 software used by Ceccato et al. is **Disco**: this is commercial soft-
1340 ware, with an academic license lasting only 30 days (although free of
1341 charge), released for Windows and MacOS operating systems only,
1342 which makes its use under Linux systems impossible except by us-
1343 ing emulation environments such as Wine.

1344 For what is my personal vision and training as a computer scientist,
1345 it would have been preferable to use a *cross-platform, freely licensed*
1346 *open source* software alternative to Disco: one such software could
1347 have been **ProM Tools** [58], a framework for process mining very
1348 similar to Disco in functionality, but fitting the criteria just described,
1349 or use Python libraries such as **PM4PY** [59], which offer ready-to-use
1350 algorithms suitable for various process mining needs.

58 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

1351 **Invariants Inference and Analysis** The limitation in this case is principally Daikon: this software is designed to compute the invariants of a software from its live execution or from a file containing its execution flow, not to find the invariants of a cyber physical system. Since there are currently no better consolidated alternatives for inferring invariants, however, an attempt was still made to use Daikon as best as possible.

```
daikon_results_cond.txt
~/UniVr/Tesi/PLC-RE/daikon/Daikon_Invariants

daikon version 5.8.14, released October 6, 2022; http://plse.cs.washington.edu/daikon.
Reading splitter info files
(read 1 spinfo file, 1 splitter)
Reading declaration files .aprogram.point:::POINT: 1 of 1 splitters successful

(read 1 decls file)
Processing trace data; reading 1 dtrace file:

Warning: No non-obvious non-suppressed exclusive invariants found in
aprogram.point:::POINT
Warning: No non-obvious non-suppressed exclusive invariants found in
aprogram.point:::POINT
=====
aprogram.point:::POINT
PLC2\_MemoryRegisters\_MW1 == PLC3\_MemoryRegisters\_MW1
PLC1\_MemoryRegisters\_MW0 == 40.0
PLC1\_MemoryRegisters\_MW1 == 80.0
PLC1\_Coils\_QX00 one of { 0.0, 1.0 }
PLC1\_Coils\_QX01 one of { 0.0, 1.0 }
PLC1\_Coils\_QX02 one of { 0.0, 1.0 }
PLC2\_MemoryRegisters\_MW1 == 10.0
PLC2\_MemoryRegisters\_MW2 == 20.0
PLC2\_Coils\_QX00 one of { 0.0, 1.0 }
PLC3\_InputRegisters\_IW0 >= 0.0
PLC3\_Coils\_QX00 one of { 0.0, 1.0 }
PLC3\_Coils\_QX02 one of { 0.0, 1.0 }
prev\_PLC1\_Coils\_QX00 one of { 0.0, 1.0 }
prev\_PLC1\_Coils\_QX01 one of { 0.0, 1.0 }
prev\_PLC2\_Coils\_QX00 one of { 0.0, 1.0 }
prev\_PLC3\_InputRegisters\_IW0 >= 0.0
prev\_PLC3\_Coils\_QX00 one of { 0.0, 1.0 }
prev\_PLC3\_Coils\_QX02 one of { 0.0, 1.0 }
PLC1\_Max\_safety == 77.0
```

Figure 3.10: Example of Daikon's output

1357 **First limitation** The biggest problem with Daikon applied to the computation of invariants of an industrial system is the difficult reading 1358 of the resulting output: the software in fact returns a very long list 1359 of invariants, one invariant per line, many of no use and without 1360 correlating invariants that may have common features or deriving 1361

1362 additional information from them. The process of screening and rec-
1363ognizing the significant invariants, as well as the correlation between
1364them, must be done by a human: certainly not an easy task given the
1365volume of invariants one could theoretically be faced with (hundreds
1366and hundreds of invariants). An example of Daikon's output can be
1367seen in Figure 3.10.

1368 **Second limitation** The bash script used in this phase of the analysis does
1369not help at all in deriving significant invariants more easily: it merely
1370launches Daikon and saves its output to a text file by simply redirect-
1371ing the stdout to file. No data reprocessing is done during this step.
1372In addition, if a condition is to be specified to Daikon before perform-
1373ing the analysis, it is necessary each time to edit the .spinfo file by
1374manually entering the desired rule, an inconvenient operation when
1375multiple analyses are to be performed with different conditions each
1376time.

1377 Table 3.1 provides a summary of the limitations discussed regarding the
1378 Ceccato et al. framework:

Phase	Limitations
Testbed	<ul style="list-style-type: none">- Oversimplified model compared other testbeds- Physics of the liquid not considered: this causes data perturbation.
Pre-processing	<ul style="list-style-type: none">- It is not possible to select a subsystem by (groups of) PLCs or by time range- Some additional attributes are manually inserted into dataset.
Graphical/Statistical Analysis	<ul style="list-style-type: none">- Only one chart at the time is displayed: difficulty in capturing the relationship between actuators and sensors.- Statistical Analysis provides little information

60 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

Business Process Analysis	- Ad hoc solution designed to work under special conditions - Use of commercial software for process mining
Invariant Analysis	- Reading output is challenging - Script for analysis merely launches Daikon without reprocessing outcomes

Table 3.1: Summary table of Ceccato et al. framework limitations

Extending and Generalizing Ceccato et al.'s Framework

¹³⁷⁹ **I**N Chapter 3, we presented the state of the art of *process comprehension*
¹³⁸⁰ of an Industrial Control System (ICS) with a focus on the methodology
¹³⁸¹ proposed by Ceccato et al. [12][Section 3.2], explaining what it consists of,
¹³⁸² its practical application on a testbed, and most importantly highlighting
¹³⁸³ its limitations and critical issues (see Section 3.2.7).

¹³⁸⁴ In this chapter we will present a **proposal to improve the methodology**
¹³⁸⁵ seen in the previous chapter, addressing most of the critical issues (or
¹³⁸⁶ at least trying to do so) described in Section 3.2.7 by almost completely
¹³⁸⁷ rewriting the original framework, enhancing its functionalities and in-
¹³⁸⁸ serting new ones where possible, while preserving its general structure
¹³⁸⁹ and approach. Indeed, the system analysis will encompass the same **four**
¹³⁹⁰ **phases** as the original methodology (Data Pre-processing, Graphs and Sta-
¹³⁹¹ tistical Analysis, Invariant Inference and Analysis, and Business Process
¹³⁹² Mining and Analysis). In addition to these phases, a **fifth phase** dedicated
¹³⁹³ exclusively to **network traffic analysis** will be introduced. Each phase of
¹³⁹⁴ the original methodology will undergo thorough revision to enhance the
¹³⁹⁵ understanding of the industrial system being analyzed and its behavior,
¹³⁹⁶ aiming to provide a more complete and coherent process comprehension.
¹³⁹⁷ This revision aims to enrich the analysis process, making it more robust

¹³⁹⁸ and transparent.

¹³⁹⁹ Please note that our proposals do not include improvements to the data
¹⁴⁰⁰ gathering phase. This decision is solely because the new framework will
¹⁴⁰¹ not be tested on the same case study used by Ceccato et al. (Section 3.2.1),
¹⁴⁰² but rather on a different case study known as the iTrust SWaT system [15].
¹⁴⁰³ iTrust has already provided some datasets that contain the execution trace
¹⁴⁰⁴ of the physical system and the network traffic scan for this case study. In
¹⁴⁰⁵ contrast to the case study conducted by Ceccato et al., the iTrust SWaT
¹⁴⁰⁶ system is not simulated but rather a **real-world system**. This distinction
¹⁴⁰⁷ is important as it introduces additional complexities and considerations
¹⁴⁰⁸ when analyzing and interpreting the data.

¹⁴⁰⁹ **Outline** The upcoming sections provide an outline of what to expect:

- ¹⁴¹⁰ • **Section 4.1** will introduce a **novel framework** that we have devel-
¹⁴¹¹ oped to address the limitations of Ceccato et al.'s framework. Sec-
¹⁴¹² tion will provide a brief overview of the framework's features and
¹⁴¹³ structure;
- ¹⁴¹⁴ • in **Section 4.2** the **features of our framework** will be demonstrated
¹⁴¹⁵ through their application to the different steps of the methodology.
¹⁴¹⁶ Practical examples will be used to illustrate the functionality of the
¹⁴¹⁷ framework. To facilitate this, we will utilize a more complex and
¹⁴¹⁸ detailed testbed compared to the one employed by Ceccato et al.

¹⁴¹⁹ 4.1 The Proposed New Framework

¹⁴²⁰ In our version of the framework we decided to follow a few design
¹⁴²¹ choices:

- ¹⁴²² 1. it must be implemented in a **single programming language**;
- ¹⁴²³ 2. it must be **as independent as possible of the system** to be analyzed;

1424 3. It must provide greater **flexibility and ease of use** at every stage.

1425 In the following, we discuss these three features in more detail.

1426 **Single Programming Language** The implementation of Ceccato et al.'s tool
1427 involved the use of different programming languages for each of the
1428 phases, ranging from Python to Java, and even including Bash script-
1429 ing.

1430 However, we believe that this heterogeneity introduces challenges
1431 in terms of user-friendliness and intuitiveness. It becomes more dif-
1432 ficult for users to navigate and operate the tool effectively. Fur-
1433 thermore, the utilization of multiple technologies complicates code
1434 maintenance and the addition of new features, especially when man-
1435 aged by a single person who may have expertise in only one lan-
1436 guage while being less familiar with others.

1437 Considering these factors, we have made the decision to adopt a
1438 single programming language for the framework to ensure homo-
1439 geneity, simplify usability, and facilitate code maintenance for fu-
1440 ture users. Python has been chosen as the language of choice due
1441 to its simplicity, readability, versatility, and powerfulness. Addition-
1442 ally, Python benefits from a vast ecosystem of libraries and packages
1443 that cater to various requirements, making it a suitable choice for our
1444 framework.

1445 **System Independence** One of the significant limitations we identified in
1446 Ceccato et al.'s tool, as highlighted in Section 3.2.7, is its **strong de-**
1447 **pendence on the specific testbed** being used. This means that the
1448 tool lacks the flexibility to configure parameters for analyzing differ-
1449 ent industrial systems.

1450 To address this limitation and achieve system independence, we have
1451 made a crucial improvement in our framework. We have introduced
1452 a comprehensive configuration file called *config.ini* that allows users
1453 to customize all the necessary parameters for analyzing the targeted

1454 system. This general configuration file eliminates any references to
1455 hardcoded variables or values found in the Ceccato et al.'s tool, pro-
1456 viding users with the flexibility to tailor the analysis according to
1457 their specific requirements.

1458 **Flexibility and Ease on Use** The lack of flexibility and ease of use in a tool
1459 can be a significant disadvantage, as it hampers its effectiveness and
1460 makes it difficult for users to achieve their desired outcomes. Cec-
1461 cato et al.'s tool was affected by these limitations, as it required users
1462 to run scripts from the command line without sufficient options or
1463 parameters for customizing the analysis. Consequently, the tool fell
1464 short in terms of user-friendliness and failed to provide the neces-
1465 sary flexibility to accommodate specific user requirements.

1466 To settle these issues, we enhanced the command-line interface in the
1467 proposed framework by adding new options and parameters. These
1468 new features provide the user with greater flexibility, enabling to
1469 specify parameters and options that allow for more in-depth anal-
1470 ysis and focused results analyzing data more effectively and effi-
1471 ciently. With these enhancements, the framework has become more
1472 user-friendly, reducing the learning curve and making it more acces-
1473 sible to a wider range of users.

1474 This, in turn, makes the framework more valuable and useful, in-
1475 creasing its adoption and effectiveness across a range of industrial
1476 control systems and applications.

1477 Moreover, with new options and parameters users can now access a
1478 range of customizable options and parameters, making the tool more
1479 intuitive and user-friendly.

1480 Overall, the enhancements made to the framework represent a signifi-
1481 cant step forward in making it more effective, efficient, and user-friendly.

4.1.1 Framework Structure

The proposed framework follows a similar structure to the original tool, with a division into five main directories representing different phases of the analysis: **data pre-processing**, **graphs and statistical analysis**, **invariant analysis**, and **process mining**. A new phase is added compared to the original, concerning the **analysis of the network traffic**. These directories contain the corresponding Python scripts responsible for performing the analysis, along with any necessary subdirectories and input/output files to ensure the proper functioning of the framework.

```
1491 .
1492   |-- config.ini
1493   |-- daikon
1494     |-- Daikon_Invariants
1495     |-- daikonAnalysis.py
1496     |-- runDaikon.py
1497   |-- network-analysis
1498     |-- data
1499     |-- networkAnalysis.py
1500     |-- export_pcap_data.py
1501     |-- swat_csv_extractor.py
1502   |-- pre-processing
1503     |-- mergeDatasets.py
1504     |-- system_info.py
1505   |-- process-mining
1506     |-- data
1507     |-- process_mining.py
1508   |-- statistical-graphs
1509     |-- histPlots_Stats.py
1510     |-- runChartSubPlots.py
```

Listing 4.1: Novel Framework structure and Python scripts

The *config.ini* file is located in the root directory of the framework. This file holds significant importance as it grants the framework independence from the industrial control system being analyzed. Within this file, users have the opportunity to configure various general parameters and options. These include specifying file paths for reading or writing files, as

1516 well as options related to specific analysis phases.

1517 The file is organized into sections, with each section dedicated to a spe-
1518 cific aspect of the configuration. These sections contain user-customizable
1519 parameters that are later referenced within the Python scripts comprising
1520 the framework. Sections of *config.ini* are:

- 1521 • [PATHS]: defines general paths such as the project root directory;
- 1522 • [PREPROC]: includes parameters needed for the **pre-processing phase**,
1523 like the directory containing the raw datasets of the individual PLCs
1524 and *granularity* for the slope calculation. The granulariy is given in
1525 terms of the time interval over which the slope is calculated;
- 1526 • [DATASET]: defines settings and parameters used during the **dataset**
1527 **enrichment** stage, for example the additional attributes;
- 1528 • [DAIKON]: defines parameters needed for **invariant analysis** with
1529 Daikon, e.g. directories and files containing the outcomes of the anal-
1530 ysis;
- 1531 • [MINING]: contains parameters used during the **process mining**
1532 phase, such as data directory;
- 1533 • [NETWORK]: includes specific settings for extracting the data ob-
1534 tained from the packet sniffing phase on the ICS network and con-
1535 verting it to CSV format. It also defines the **network protocols** that
1536 are to be analyzed.

1537 4.1.2 Python Libraries and External Tools

1538 As the framework has been developed entirely in Python, the objec-
1539 tive was to minimize reliance on external tools and instead integrate vari-
1540 ous functionalities within the framework itself. The aim was to make the
1541 framework independent from external software. The only remaining ex-
1542 ternal tool from the Ceccato et al. tool is Daikon. This choice was made

1543 because there is currently no better alternative or Python package avail-
1544 able that offers the same functionalities as Daikon.

1545 Instead, the framework extensively utilizes Python libraries for han-
1546 dling various functionalities and input data. The core libraries on which
1547 the framework relies are:

- 1548 • **Pandas**, also used in the Ceccato et al.'s tool for dataset management,
1549 but whose use here has been deepened and extended
- 1550 • **NumPy**, often used together with Pandas to perform some opera-
1551 tions to support it;
- 1552 • **MatPlotLib**, for managing and plotting graphical analysis;
- 1553 • other scientific libraries such as **SciPy**, **StatsModel** [60] and **Net-**
1554 **workX** [61], for mathematical, statistical and analysis operations on
1555 the data;
- 1556 • **GraphViz**, for the creation of activity diagrams in the process mining
1557 phase.

1558 Having now seen the structure of the framework, in the next sections we
1559 will go into more detail describing our proposal.

1560 4.2 Analysis Phases

1561 In this section, the behavior of the proposed framework throughout the
1562 various analysis phases of the methodology will be illustrated. Here is a
1563 brief **outline** of the content covered in this section:

- 1564 • **Section 4.2.1:** this section will present the **testbed** used to demon-
1565 strate examples of the framework's application;
- 1566 • **Section 4.2.2:** the introduction of the new **network traffic analysis**
1567 phase (*Phase 0*) will be discussed. This phase aims to understand the
1568 structure of the industrial system's network and analyze its commu-
1569 nication patterns;

- **Section 4.2.3:** the **pre-processing** phase (*Phase 1*) will be presented, consisting of two parts: dataset merging and enrichment, followed by a preliminary analysis of the resulting dataset;
 - **Section 4.2.4:** this section focuses on the **Graphs and Statistical Analysis** phase (*Phase 2*). It will demonstrate the extraction of valuable information from the system by analyzing graphs that represent the behavior of registers over time;
 - **Section 4.2.5:** the **Invariant Inference** phase (*Phase 3*) will be explored. This phase involves deriving system information based on discovered invariants through the use of Daikon. Two examples of semi-automatic analysis will be showcased;
 - **Section 4.2.6:** the **Business Process Mining** phase (*Phase 4*) will be covered. This phase aims to gain an overview of the system's behavior and extract additional information through process mining techniques.

1585 4.2.1 A Little Testbed: Stage 1 of iTrust SWaT System

1586 Before we proceed with presenting the analysis steps of the proposed
1587 framework, let us introduce the testbed that will serve as an illustration for
1588 practical examples, demonstrating the effectiveness of our methodology
1589 and the potential of the framework. This testbed corresponds to Stage 1 of
1590 the iTrust SWaT (Secure Water Treatment) system [15]. The selection of this
1591 testbed is intentional, as it serves as a precursor to the comprehensive case
1592 study we will be addressing in the upcoming chapters. The iTrust SWaT
1593 system offers elements of greater complexity compared to the individual
1594 stages of the Ceccato et al. testbed.

1595 The testbed comprises several components, including:

- 1596 • a **tank**, which serves as the main element of interest;

- 1597 • a PLC responsible for monitoring and controlling the operations within
1598 the stage;
- 1599 • **two sensors** that provide readings of the water level within the tank
1600 and the incoming flow. These sensors are identified as LIT101 and
1601 FIT101 in the PLC registers;
- 1602 • **three actuators**, namely a valve and two pumps. These actuators
1603 regulate the level within the tank by controlling the inflow and out-
1604 flow of the liquid. These sensors are identified as MV101 (valve), P101
1605 and P101 (pumps) in the PLC registers.

1606 Despite its moderate complexity, this testbed provides an ideal plat-
1607 form for presenting straightforward and concise examples of the frame-
1608 work's behavior. It enables us to effectively demonstrate the potential
1609 of the framework and facilitate a deeper understanding of its underlying
1610 methodology.

1611 4.2.2 Phase 0: Network Analysis

1612 The objective of the network analysis presented in this section is to offer
1613 users valuable information regarding the communication process within
1614 an industrial control system and a broader perspective on network com-
1615 munications, delving into previously unexplored aspects of the system.
1616 These additional dimensions provide a deeper understanding of the sys-
1617 tem's behavior and characteristics. This analysis aims to provide users
1618 with an overview of the communication between PLCs at level 1 (see Fig-
1619 ure 2.1), as well as the communication between PLCs and devices at higher
1620 levels such as HMIs and Historian servers (see Section 2.1 for ICSs archi-
1621 tecture). This also allows us to have a better understanding of the system
1622 architecture and network topology. The analysis focuses on industrial pro-
1623 tocols used and the information exchanged.

1624 By reconstructing the network communication structure using data ob-
1625 tained from the network traffic sniffing process, users can gain a compre-
1626 hensive understanding of the behavior of the underlying industrial sys-

1627 tem. This knowledge can then be utilized to **plan a strategy** for analyzing
1628 the physical processes within the system.

1629 **4.2.2.1 Extracting Data from PCAP Files**

1630 The initial step involves extracting the desired information from the
1631 PCAP files that contain the captured network traffic. This includes details
1632 such as the source IP address, destination IP address, protocol used, and
1633 the type of request made (e.g., Read/Write, Request/Response). The ex-
1634 tracted data is then converted into a more convenient CSV format. This
1635 extracted data serves as the foundation for the subsequent phase.

1636 In the latter part of this phase, the extracted data is utilized to generate
1637 the network schema. The network schema provides a visual representa-
1638 tion of the connections and relationships within the network, showcasing
1639 the communication patterns between different components. This schema
1640 helps in understanding the overall structure and behavior of the industrial
1641 control system.

1642 To accomplish the extraction of data from the PCAP files, a Python
1643 script called `export_pcap_data.py` is employed. This script, originally de-
1644 signed for the business process phase, is located in the directory
1645 `$(project_dir)/network-analysis` and accepts the following options as
1646 command-line arguments:

- 1647 • **-f or --filename:** allows the user to specify a single PCAP file to be
1648 passed as input to the script. The user can provide the complete file
1649 path of the PCAP file as an argument;
- 1650 • **-m or --mergefiles:** enables the merging of multiple PCAP files. In
1651 this scenario, the files should be located within the directory speci-
1652 fied by the `pcap_dir` directive in the `config.ini` configuration file and
1653 the user does not have to provide the path to each PCAP file;
- 1654 • **-d or --mergedir:** allows for specifying the directory that contains the
1655 PCAP files to be automatically imported into the script and merged.

1656 This ensures that all the PCAP files within the specified directory will
1657 be processed by the script without the need for manual selection or
1658 input.

- 1659 • **-s or --singledir:** operates differently from the previous option men-
1660 tioned. This option enables the extraction of data from each individ-
1661 ual PCAP file within the specified directory. The extracted data is
1662 then saved in separate CSV datasets, which are stored in the direc-
1663 tory specified by the `split_dir` directive in the `config.ini` file. This
1664 functionality proves useful when dealing with exceptionally large
1665 PCAP files, where merging them together for export might consume
1666 significant time and resources. By utilizing this option, the extrac-
1667 tion procedure becomes lighter and more manageable. The extracted
1668 data in separate CSV files can be utilized in the later stages of the
1669 Network Analysis process;
- 1670 • **-t or --timerange:** this functionality enables users to specify a specific
1671 time period within the PCAP files from which they wish to extract
1672 relevant information.

1673 Unless the `-s` option is explicitly specified, the results of data extraction
1674 and export will be saved to a single CSV dataset within the
1675 `$(project_dir)/network-analysis/data` directory. The default file name
1676 for this output file is determined by the `pcap_export_output` directive spec-
1677 ified in the `config.ini` file. In addition, by utilizing the `protocols` and
1678 `ws_<protocol>_field` directives, user can configure the network protocols
1679 to be searched within the PCAP files. Furthermore, user can specify the
1680 relevant Tshark/Wireshark fields to extract for the specified protocols set
1681 in the `protocols` directive.

1682 After obtaining the extracted data, it is possible to proceed with the
1683 second part of the network analysis.

1684 **4.2.2.2 Network Information**

1685 During this stage, the exported CSV data is processed to derive valua-
1686 able information regarding network communications and the structure of
1687 the network itself. The objective is to identify and establish relationships
1688 between IP addresses present on the network, thereby determining the
1689 sources and destinations of communications. Furthermore, the analysis
1690 detects the protocols used for each communication and quantifies the var-
1691 ious types of requests made.

1692 This information is then transformed into a **graph representation of**
1693 **the network** (or subnetwork, if specified). In this graph, devices are repre-
1694 sented as nodes labeled with their IP addresses, while edges represent the
1695 incoming and outgoing communications of these devices, along with the
1696 corresponding information.

1697 To ensure comprehensibility, the analysis also provides users with **tex-**
1698 **tual information** containing the same details as the graph representation.
1699 This text-based information serves as an alternative for cases where the
1700 graph may become complex to interpret, particularly when numerous edges
1701 connect nodes and result in a high volume of network requests.
1702 This textual information is saved to another CSV file, enabling offline ref-
1703 erence or potential future utilization. By having this file available, users
1704 can access the network analysis results in a structured format for further
1705 analysis or documentation purposes.

1706 The Python script `networkAnalysis.py` in the `$(project_dir)/network-analysis`
1707 directory manages this phase of the analysis. The script can be executed
1708 with the following parameters:

- 1709 • **-f** or **--filename**: used to specify the CSV dataset containing the net-
1710 work data exported in the previous step. The dataset should be lo-
1711 cated in the directory `$(project_dir)/network-analysis/data`;
- 1712 • **-D** or **--directory**: used to specify the directory that contains the CSV
1713 datasets obtained using the **-s** option of the Python script `export_pcap_data.py`.

1714 By passing this parameter, the script will automatically merge the
1715 datasets and proceed with the analysis of the data contained within
1716 them;

- 1717 • **-s or --srcaddr:** allows for specifying the source IP address for which
1718 you wish to display the incoming and outgoing communications. By
1719 providing the source IP address as an argument, the script will focus
1720 on showcasing the communications associated with that particular
1721 IP address;
- 1722 • **-d or --dstaddr:** enables the user to specify the destination IP address
1723 for which you want to display the incoming and outgoing communi-
1724 cations. By providing the destination IP address as an argument, the
1725 script will concentrate on presenting the communications associated
1726 with that specific IP address.

1727 The parameters related to IP addresses, including source and destina-
1728 tion, are optional. It is possible to specify either one of them individually.
1729 For instance, if the user specifies only the source IP address, the script will
1730 display the network nodes with which it communicates on the outgoing
1731 side, along with the corresponding generated traffic. Similarly, if only the
1732 destination IP address is specified, the script will showcase the network
1733 nodes communicating with it on the incoming side, along with the rele-
1734 vant traffic data.

1735 During the analysis, the script identifies and displays the IP addresses
1736 present in the network as output for the user's reference. This allows the
1737 user to select specific IP addresses from the command line for a more fo-
1738 cused analysis, such as choosing a subnet of interest. Additionally, the
1739 script detects and tracks distinct communications between pairs of PLCs,
1740 keeping a record of the number of these communications.

1741 The result of the analysis is a graph representation of the network (or
1742 subnetwork) to be analyzed. An example of such a graph can be seen in
1743 Figure 4.1.

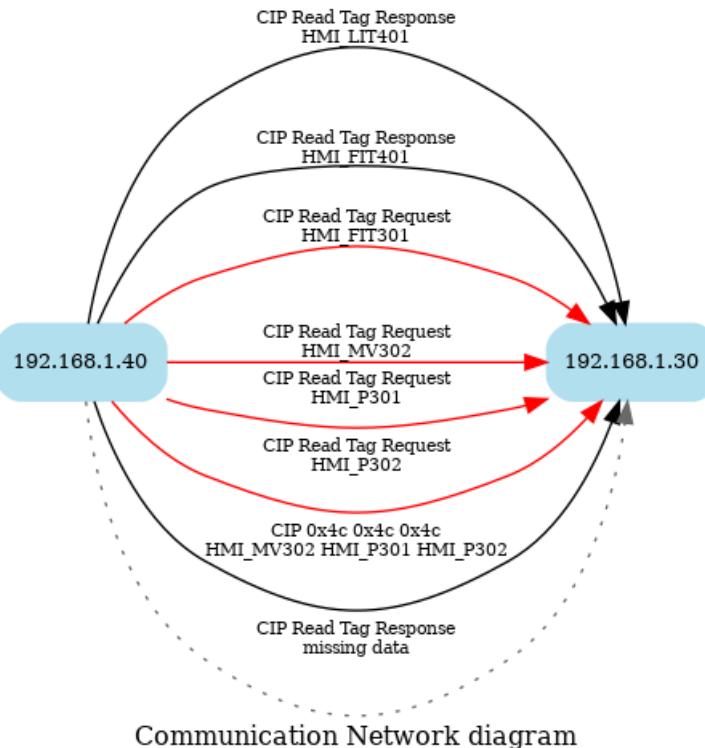


Figure 4.1: Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination)

1744 The graph illustrates the network communications between the source
 1745 IP address 192.168.1.40 and the destination IP address 192.168.1.30. Each
 1746 arrow represents a communication between these IP addresses, showcas-
 1747 ing the flow and direction of the interactions. The red arrows indicate re-
 1748 quests initiated by the source IP address towards the destination IP, while
 1749 the black arrows represent responses sent by the source IP in response
 1750 to previous requests made by the destination IP. The gray dotted arrows
 1751 represent responses for which the corresponding request is missing or un-
 1752 available for some reason. Overall, the graph distinguishes the different
 1753 types of communications and provides insights into the request-response
 1754 dynamics between the source and destination IP addresses. The graph is
 1755 automatically generated and saved within the
 1756 `$(project_dir)/network-analysis/data` directory.

1757 In terms of the textual output, we can observe how the same data is
 1758 represented. The communications exchanged between the two PLCs are
 1759 displayed more prominently, allowing for a clearer understanding. Unlike
 1760 the graph, the textual representation includes a column on the right-hand
 1761 side, indicating the number of communications for each type of request.
 1762 This provides a more distinct perspective on the network behavior within
 1763 an industrial control system that utilizes, in this case, the CIP protocol for
 1764 its communications.

src	dst	protocol	service_detail	register	
192.168.1.40	192.168.1.30	CIP	Read Tag Response	HMI_LIT401	11249
				HMI_FIT401	10539
			Read Tag Request	HMI_FIT301	8031
				HMI_MV302	7209
				HMI_P301	7115
				HMI_P302	7040
			0x4c 0x4c 0x4c	HMI_MV302 HMI_P301 HMI_P302	1
			Read Tag Response	missing data	1

Figure 4.2: Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination) in textual mode

1765 As previously mentioned, the textual data is stored in a CSV dataset
 1766 located in the \$(project_dir)/network-analysis/data directory.

1767 4.2.3 Phase 1: Data Pre-processing

1768 *Data Pre-processing phase* is probably the most delicate and significant
 1769 one: depending on how large the industrial system to be analyzed is, the
 1770 data collected, and how it is enriched using the additional attributes, the
 1771 subsequent system analysis will provide more or less accurate outcomes.

1772 The Ceccato et al.'s tool has several limitations, particularly at this
 1773 stage. It does not allow for the isolation of a subsystem, either in terms
 1774 of time or the number of PLCs to be analyzed. The system is considered as
 1775 a whole without the ability to focus on specific subsystems. Additionally,
 1776 many of the additional attributes had to be manually added, and for the
 1777 ones entered automatically, there is no way to specify the register type to
 1778 associate them with.

1779 The combination of these limitations, along with the presence of hard-
1780 coded references to attributes and registers in the tool’s code, makes the
1781 analysis of the system more challenging. Furthermore, it compromises the
1782 accuracy and reliability of the obtained results in terms of both quantity
1783 and quality.

1784 In the proposed framework, these issues have been addressed by in-
1785 corporating new features.

1786 *Firstly*, the framework allows for the **selection of a subsystem from**
1787 **the command line** based on both temporal criteria and the specific PLCs
1788 to be included. This enables more focused and targeted analysis.

1789 *Secondly*, we have **revamped the process of enriching** the resulting
1790 dataset by eliminating manual entry of additional attributes. Instead, users
1791 now have the flexibility to determine the type of additional attribute to as-
1792 sociate with a specific register.

1793 *Thirdly*, after the pre-processing stage, a **preliminary analysis** can be
1794 conducted on the resulting dataset. This analysis aims to identify the reg-
1795 isters that are associated with actuators, measurements, and hardcoded
1796 setpoints or constants. It provides insights into the dataset and helps in
1797 refining the enrichment step. The parameters for this analysis can be con-
1798 figured in the *config.ini* file, allowing for customization and fine-tuning of
1799 the process.

1800 In the upcoming sections, we will delve into a more comprehensive
1801 examination of the achievements made in this phase.

1802 4.2.3.1 Subsystem Selection

1803 In Ceccato et al.’s tool, the datasets for each individual PLC in CSV for-
1804 mat were required to be placed in a specific directory that was hardcoded
1805 in the script. The script would then merge and enrich these datasets to
1806 generate a single output dataset representing the complete process trace of
1807 the industrial system. However, the script did not provide options to se-
1808 lect specific PLCs for analysis or define a temporal range for analysis. This

lack of flexibility made the analysis more complex, especially when dealing with *transient states* (i.e., general states in which the industrial system is still initializing before actually reaching full operation) or when focusing on specific parts of the industrial system during certain periods of interest. The fixed dataset structure also may increase the number of variables that could be analyzed.

Furthermore, the tool in question did *not* allow for specifying an output CSV file to save the resulting dataset. Each dataset creation and enrichment operation would overwrite the previous file, making it inconvenient for comparisons between different execution traces unless the files were manually renamed.

The proposed framework addresses these issues by introducing improvements. First of all, in the general *config.ini* file there are some general default settings about paths, and among them the one concerning the directory where to place the datasets to be processed. In addition to this option, there are other ones that define further aspects related to the operations performed in this phase. Listing 4.2 shows the settings in question:

```
[PATHS]
root_dir = /home/marcuzzo/UniVr/Tesi
project_dir = %(root_dir)s/PLC-RE
net_csv_path = %(root_dir)s/datasets_SWaT/2015/Network_CSV

[PREPROC]
raw_dataset_directory = datasets_SWaT/2015 # Directory
    ↳ containing datasets
dataset_file = PLC_SWaT_Dataset.csv # Default output
    ↳ dataset
granularity = 10 # slope granularity
number_of_rows = 20000 # Seconds to consider
skip_rows = 100000 # Skip seconds from beginning
```

Listing 4.2: Paths and parameters for the Pre-processing phase in config.ini file

At the same time, the user has the option to specify these settings via the command line using the new Python script called *mergeDatasets.py*, located in the pre-processing directory of the project. Any options provided

1842 through the command line will override the default settings specified in
1843 the *config.ini* file. These options are:

- 1844 • **-s or --skiprows:** initial transient period (expressed in seconds) to
1845 be skipped. This option is useful in case the system has an initial
1846 transient or the analyzer wishes to start the analysis from a specific
1847 point in the dataset;
- 1848 • **-n or --nrows:** time interval under analysis, expressed in terms of the
1849 number of rows in the dataset.
1850 This option makes a **selection** on the data of the dataset;
- 1851 • **-p or --plcs:** PLCs to be merged and enriched. The user can specify
1852 the desired PLCs by indicating the CSV file names of the associated
1853 datasets with no limitations on number.
1854 This option makes a **projection** on the data of the dataset.
- 1855 • **-d or --directory:** performs the merge and enrichment of all CSV files
1856 contained in the directory specified by user, overriding the default
1857 setting in *config.ini*. It is in fact the old functionality of the Ceccato
1858 et al.'s tool, maintained here to give the user more flexibility and
1859 convenience in case he wants to perform the analysis on the whole
1860 system. This is also the default behavior in case the -p option is not
1861 specified.
- 1862 • **-o or --output:** specifies the name of the file in which the obtained
1863 dataset will be saved. It must necessarily be a file in CSV format.
- 1864 • **-g or --granularity:** specifies a granularity (expressed in seconds)
1865 that will be used to calculate the measurement slope during the dataset
1866 enrichment phase. We will discuss this later in Section 4.2.3.2.

1867 4.2.3.2 Dataset Enrichment

1868 After a step in which a function is applied to each PLC-related dataset
1869 to eliminate its unused registers within the system¹, the **dataset enrichment**
1870 **operation** is performed.

1871 This operation brings about several distinctions compared to the pre-
1872 vious version. Firstly, it is performed on each individual dataset instead
1873 of the resulting dataset. Additionally, there are a greater number of at-
1874 tributes, which are automatically calculated and added to the dataset by
1875 the `mergeDatasets.py` script. Importantly, the configuration file `config.ini`
1876 under the [DATASET] section allows users to determine which registers
1877 should be assigned to these attributes.

1878 In Listing 4.3 we can see the list of additional attributes and how they
1879 should be associated with the registers of the dataset:

```
1880 [DATASET]
1881 timestamp_col = Timestamp
1882 max_prefix = max_
1883 min_prefix = min_
1884 max_min_cols_list = lit|ait|dpit
1885 prev_cols_prefix = prev_
1886 prev_cols_list = mv[0-9]{3}|p[0-9]{3}
1887 trend_cols_prefix = trend_
1888 trend_cols_list = lit
1889 trend_period = 150
1890 slope_cols_prefix = slope_
1891 slope_cols_list = lit
```

Listing 4.3: config.ini parameters for dataset enriching

1892 In the following, we report a brief explanation of the parameters just seen:

1893 **timestamp_col** denotes the name of the column in the dataset that holds
1894 the timestamp data. This parameter is significant not only in the
1895 current phase but also in the Process Mining phase. In the Ceccato et

¹This becomes particularly relevant when conducting a Modbus register scan, where ranges of registers are examined. In this process, it is assumed that any unused registers hold a constant value of zero.

1896 al.'s work, this parameter was hardcoded and lacked configurability,
1897 leading to errors if the analyzed system changed.

1898 **max_prefix, min_prefix, max_min_cols_list** refer to any relative maximum
1899 or minimum values (*relative setpoints*) of one or more measures and
1900 that can be found and inserted as new columns within the dataset.
1901 The first two parameters indicate the prefix to be used in the column
1902 names affected by this additional attribute, while the third specifies
1903 of which type of registers we want to know the maximum and/or
1904 minimum value reached (several options can be specified using the
1905 logical operator | - or).

1906 If, for example, we want to know the maximum value of the registers
1907 associated with the tanks, indicated in the iTrust SWaT system by the
1908 prefix LIT, we only need to specify the necessary parameter in the
1909 *config.ini* file, so `max_min_cols_list = lit`.

1910 The result will be to have in the dataset thus enriched a new column
1911 named `max_LIT101`.

1912 **prev_cols_prefix, prev_cols_list** refer to the values at the previous time
1913 instant of the registers specified in `prev_cols_list`. It is possible to
1914 specify registers using *regex*, as in the example shown. It may be use-
1915 ful in some cases to have this value available to check, for example,
1916 when a change of state of a single given actuator occurs. The behav-
1917 ior of these parameters is the same as described in the point above.

1918 **slope_cols_prefix, slope_cols_list** are related to the calculation of the
1919 slope of a specific register that contains numeric values (usually a
1920 measure), that is, its trend. Slope calculation makes little sense on
1921 booleans. The slope can be **ascending** (if its value is greater than
1922 zero), **descending** (if less than zero) or **stable** (if approximately equal
1923 to zero). We will delve into the details of slope calculation in the fol-
1924 lowing paragraph, as it pertains to the attributes `trend_cols_prefix`,
1925 `trend_cols_list`, and `trend_period`.

Initially, the parameters for registers to be associated with each additional attribute may be left blank, as we may not have prior knowledge about the system and are unsure about which registers correspond to actuators, measurements, or other attributes. This information can be obtained from the preliminary analysis that follows the merging of datasets. The analysis, performed based on user's choice, provides indications on potential sensors, actuators, and other relevant information. These indications help the user set the desired values in the *config.ini* file and refine the enrichment process by re-launching the *mergeDatasets.py* script.

Slope Calculation The *slope* is an attribute that represents the **trend** of the measurement being considered. It is particularly useful, in our context, during the inference and invariant analysis phase to gather information about the trend under specific conditions. The slope can generally be classified as **increasing** ($\text{slope} > 0$), **decreasing** ($\text{slope} < 0$), or **stable** ($\text{slope} = 0$).

Normally, the slope is calculated through a simple mathematical formula: given an interval a, b relative to the measurement l , the slope is given by the difference of these two values divided by the amount of time t that the measurement takes to reach b from a :

$$\text{slope} = \frac{l(b) - l(a)}{t(b) - t(a)}$$

In the proposed framework, similar to the Ceccato et al.'s tool, this time interval (the granularity) can be adjusted to be either long or short.

The choice of granularity depends on the desired accuracy of the slope calculation. A lower granularity will provide a slope that closely reflects the actual measurement trend, while a higher granularity will result in flatter slope data. Each time interval within which the measurement is divided corresponds to a slope value. These slopes are calculated and added as additional attributes in the dataset. Later on, these slope values are used to determine the trend of the measurement in specific situations

¹⁹⁵⁴ or conditions.

¹⁹⁵⁵ Calculating the slope directly from the raw measurement data can be
¹⁹⁵⁶ a suitable approach for systems where the measurements are *not* heavily
¹⁹⁵⁷ influenced by **perturbations**. Perturbations, such as liquid oscillations in
¹⁹⁵⁸ a tank during filling and emptying phases, can lead to fluctuating read-
¹⁹⁵⁹ ings of the level. In such cases, maintaining a low granularity can provide
¹⁹⁶⁰ a more accurate calculation of the overall trend that closely aligns with
¹⁹⁶¹ the actual measurement trend. This situation occurs, for instance, in the
¹⁹⁶² testbed utilized by Ceccato et al.

¹⁹⁶³ However, if perturbations significantly affect the measurement read-
¹⁹⁶⁴ ings, calculating the slope on individual time intervals may result in an
¹⁹⁶⁵ inaccurate trend definition, irrespective of the chosen granularity. In such
¹⁹⁶⁶ cases, the fluctuating nature of the measurements due to perturbations can
¹⁹⁶⁷ introduce errors in the slope calculation, making it less reliable as an indi-
¹⁹⁶⁸ cator of the actual trend.

¹⁹⁶⁹ Figure 4.3 demonstrates this assertion: the measurement, in blue, refers
¹⁹⁷⁰ to the LIT101 tank of our testbed; in red, the slope calculation related to
¹⁹⁷¹ the measurement with three different granularities: 30 (Figure 4.3a), 60
¹⁹⁷² (Figure 4.3b) and 120 seconds (Figure 4.3c).

¹⁹⁷³ It is noticeable that as the granularity increases, the slope values flatten.
¹⁹⁷⁴ Moreover, in the time interval between seconds 1800 and 4200, the level of
¹⁹⁷⁵ LIT101 exhibits a predominantly increasing trend, yet the calculated slope
¹⁹⁷⁶ values fluctuate between positive and negative. Consequently, during the
¹⁹⁷⁷ invariant analysis, the overall increasing trend may not be detected, re-
¹⁹⁷⁸ sulting in a loss of information.

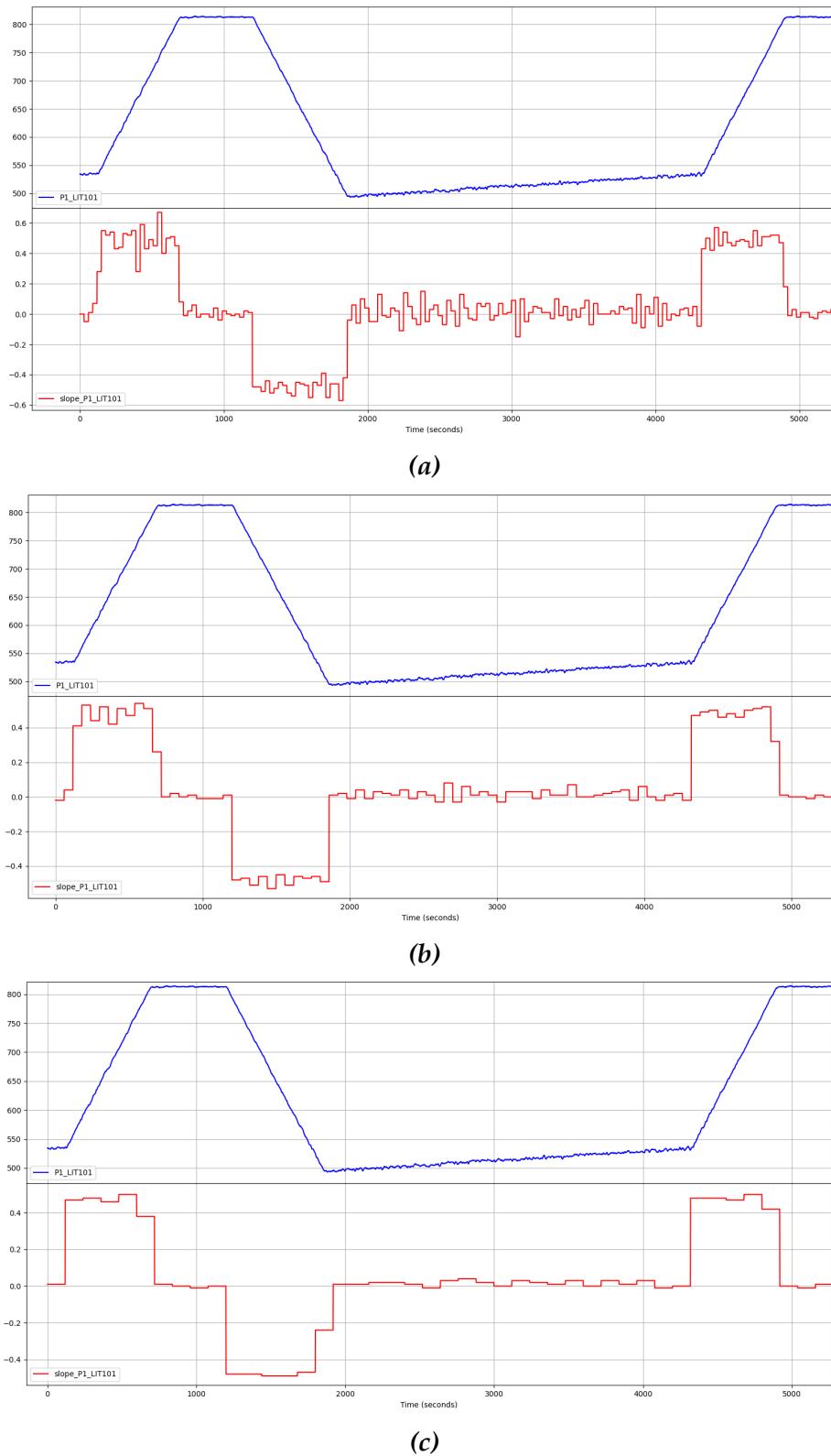


Figure 4.3: Slope comparison with granularity 30 (a), 60 (b) and 120 seconds (c)

1979 Ceccato et al.'s tool did not take into account the possibility of having
1980 strongly perturbed data, which presented a challenge that we needed to
1981 address in the development of the proposed framework.

1982 The solution to this problem involves applying techniques to reduce
1983 the "noise" in the data, aiming to achieve a more linear trend in the mea-
1984 surement curve. By minimizing the effects of perturbations, we can calcu-
1985 late slopes more accurately.

1986 There are various methods available for smoothing out noise in the
1987 data. In our framework, we focused on two commonly used approaches
1988 found in the literature: **polynomial regression** and **seasonal decomposi-**
1989 **tion**. In addition to these two methods, we also explored the use of a **line**
1990 **simplification algorithm**.

1991 *Polynomial regression* [62] is a technique that allows us to create a filter
1992 to reduce the impact of noise on the data. By fitting a polynomial function
1993 to the measurements, we can obtain a smoother curve that captures the
1994 underlying trend while minimizing the effects of perturbations.

1995 *Seasonal decomposition* [63], specifically the part related to trending, is
1996 another method we explored. It involves decomposing the time series into
1997 different components, such as trend, seasonality, and residual. By isolat-
1998 ing the trend component, we can obtain a cleaner representation of the
1999 underlying pattern in the data.

2000 *Line simplification algorithms* [64] aim to reduce the complexity of a poly-
2001 line or curve by approximating it with a simplified version composed of
2002 fewer points. By selectively removing redundant or less significant points,
2003 line simplification algorithms help reduce storage space and computa-
2004 tional requirements while preserving the overall shape and characteristics
2005 of the original line.

2006 Regarding polynomial regression, we evaluated the use of the **Savitzky-**
2007 **Golay filter** [65] as a smoothing technique. For seasonal decomposition,
2008 we explored the **Seasonal-Trend decomposition using LOESS** (STL) method
2009 [66]. For the line simplification algorithm, we specifically considered the

2010 **Ramer-Douglas-Peucker (RDP) algorithm [67].**

2011 Figure 4.4 shows a quick graphical comparison of these techniques
2012 compared with the original data. The solution adopted is the *STL decomposition*
2013 method, which effectively reduces noise compared to the Savitzky-
2014 Golay filter. However, it should be noted that this method may introduce
2015 some delay in certain parts of the data, as is typically observed in similar
2016 algorithms. Despite its apparent effectiveness, the RDP algorithm fails to
2017 accurately approximate sections where the measurement level remains rel-
2018 atively stable. Consequently, it yields incorrect slope estimations, causing
2019 a loss of valuable information about the system.

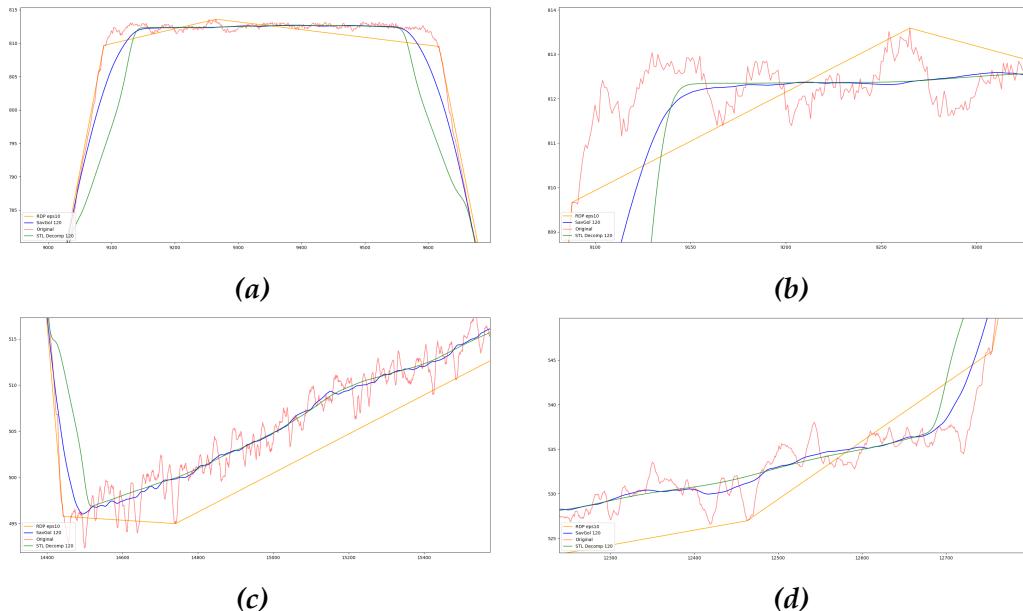


Figure 4.4: Savitzky-Golay filter (blue line), STL decomposition (green) and RDP algorithm (orange) comparison

2020 By applying the STL decomposition, we observe a notable enhance-
2021 ment in slope calculation even when using a low granularity. Figure 4.5
2022 demonstrates that, with the same granularity as shown in Figure 4.3a, the
2023 slope values, albeit exhibiting fluctuations, consistently align with the un-
2024 derlying trend of the data curve. The introduced lag resulting from the
2025 decomposition's periodicity is responsible for the observed delay.

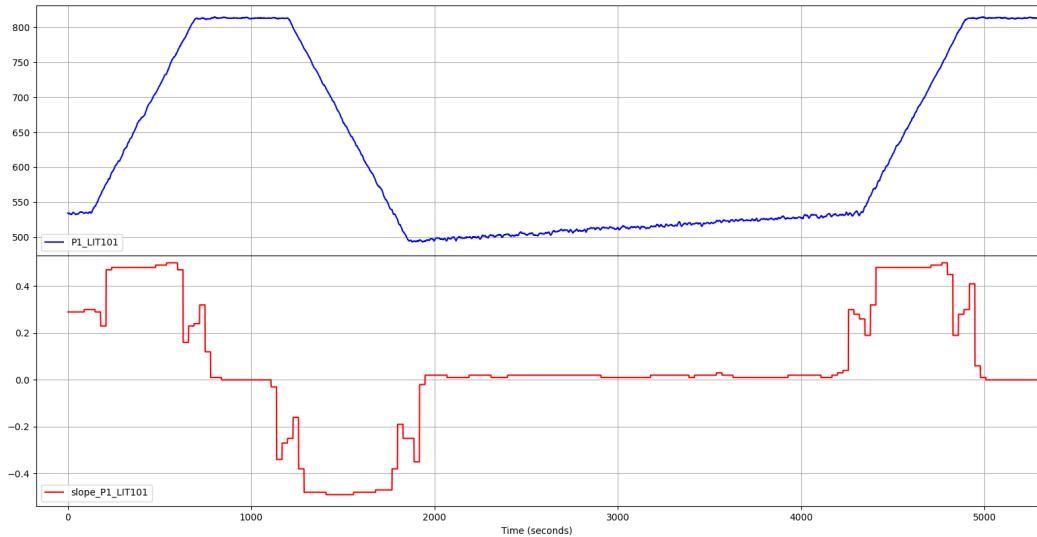


Figure 4.5: Slope after the application of the STL decomposition

2026 The periodicity, which defines the sampling time window for decom-
 2027 position and the level of noise smoothing, can be configured using the
 2028 `trend_period` directive in the `config.ini` file.
 2029 During the slope calculation, the analysis will be performed on the data
 2030 from the additional measurement trend attributes specified in the `trend_cols_list`
 2031 directive of the configuration file, rather than on the original unfiltered
 2032 data.

2033 To ensure proper interpretation by Daikon, the decimal values repre-
 2034 senting the calculated slopes are converted into **three numerical values**:
 2035 -1, 0, and 1. These values correspond to *decreasing* (if the slope is less than
 2036 zero), *stable* (if it is equal to zero), and *increasing* (if it is greater than zero)
 2037 trends, respectively. Figure 4.6 displays the modified slopes along with
 2038 the curve obtained from the STL decomposition:

2039 4.2.3.3 Datasets Merging

2040 During this step, the datasets of the individual PLCs are merged, re-
 2041 sulting in two separate datasets. The first dataset is enriched with addi-
 2042 tional attributes but excludes the timestamp column. This dataset is in-

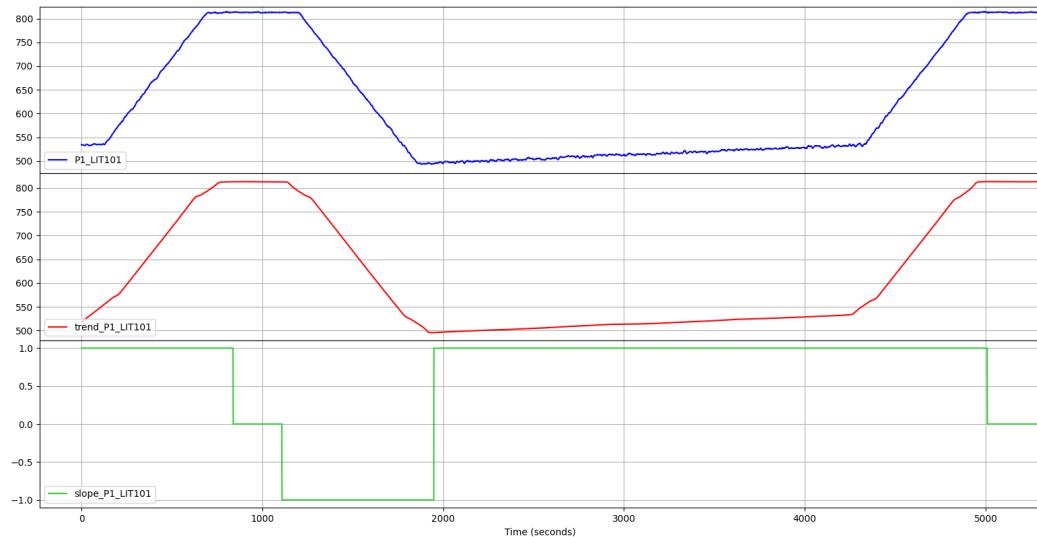


Figure 4.6: The new slope representation (green line) and the smoothed measurement data obtained with the STL decomposition (red)

2043 tended for inference and invariant analysis. The second dataset does *not*
 2044 contain any additional data and is specifically used in the process mining
 2045 phase. This dataset contains the timestamp.

2046 By default, the enriched dataset will be saved in CSV format in the
 2047 `$(project-dir)/daikon/Daikon_Invariants` directory. The other dataset,
 2048 without additional data, will be saved in the `$(project-dir)/process-mining/data`
 2049 directory. It's worth noting that both paths can be configured in the
 2050 `config.ini` file. The dataset name can be specified in the `config.ini` file or
 2051 through the `-o` command-line option. When generating the dataset for
 2052 process mining, the script will automatically add a `_TS` suffix to the file-
 2053 name to indicate that it includes the timestamp. This flexibility allows the
 2054 user to provide a different filename for each output, preventing overwrit-
 2055 ing of previous datasets. It enables the user to save the execution trace of
 2056 the selected subsystem separately and utilize them in subsequent analysis
 2057 phases.

2058 **4.2.3.4 Preliminary Analysis of the Obtained Subsystem**

2059 After merging the datasets, the user has the option to perform an **optional analysis** of the resulting dataset to extract preliminary data. This
 2060 analysis aims to gather basic information about the (sub)system and po-
 2061 tentially refine the enrichment process. If the user chooses to proceed with
 2062 the analysis, the `mergeDatasets.py` script invokes another Python script lo-
 2063 cated in the `$(project-dir)/pre-processing` directory called `system_info.py`.

2065 Relying on an analysis based on a combination of Daikon and Pandas
 2066 this script performs a quick analysis of the dataset allowing to **estimate**, al-
 2067 beit approximately, the **type of registers** (sensors, actuators, ...), also iden-
 2068 tifying possible maximum and minimum values of measurements and
 2069 hardcoded setpoints. Furthermore, leveraging the use of the additional
 2070 attribute `prev_`, the `system_info.py` script is capable of deriving measure-
 2071 ment values corresponding to state changes of individual actuators. This
 2072 allows for the identification of specific measurements associated with the
 2073 activation or deactivation of certain actuators within the system.

2074 As the last information we have duration of actuator states for each
 2075 cycle of the system: this information can be useful for making assumptions
 2076 and conjectures about the behavior of an actuator in a specific state or, by
 2077 observing the duration values of each cycle, highlighting anomalies in the
 2078 system.

2079 Listing 4.4 shows an example of the output this brief analysis related to
 2080 our testbed (for brevity, only one measurement is reported in the analysis
 2081 of actuator state changes):

```
2082     Do you want to perform a brief analysis of the dataset? [y
2083     ↵ /n]: y

2084

2085     Actuators:
2086     MV101 [0.0, 1.0, 2.0]
2087     P101 [1.0, 2.0]

2088

2089     Sensors:
2090     FIT101 {'max_lvl': 2.7, 'min_lvl': 0.0}
```

```
2091     LIT101 {'max_lvl': 815.1, 'min_lvl': 489.6}
2092
2093     Hardcoded setpoints or spare actuators:
2094     P102 [1.0]
2095
2096     Actuator state changes:
2097         LIT101      MV101      prev_MV101
2098         800.7170    0          2
2099         499.0203    0          1
2100         800.5992    0          2
2101         498.9026    0          1
2102         800.7170    0          2
2103         499.1381    0          1
2104         801.3058    0          2
2105         498.4315    0          1
2106         801.4628    0          2
2107         498.1567    0          1
2108
2109         LIT101      MV101      prev_MV101
2110         805.0741    1          0
2111         805.7414    1          0
2112         805.7806    1          0
2113         805.1133    1          0
2114         804.4068    1          0
2115
2116         LIT101      MV101      prev_MV101
2117         495.4483    2          0
2118         497.9998    2          0
2119         495.9586    2          0
2120         495.8016    2          0
2121         494.5847    2          0
2122
2123         LIT101      P101      prev_P101
2124         536.0356    1          2
2125         533.3272    1          2
2126         542.1591    1          2
2127         534.8581    1          2
2128         540.5890    1          2
2129
```

```

2130      LIT101        P101        prev_P101
2131      813.0031       2           1
2132      813.0031       2           1
2133      811.8256       2           1
2134      812.7283       2           1
2135      813.3171       2           1
2136
2137      Actuator state durations:
2138      MV101 == 0.0
2139      9   9   10   9   9   10   9   9   10   9
2140
2141      MV101 == 1.0
2142      1174  1168  1182  1160  1172
2143
2144      MV101 == 2.0
2145      669   3019  3012  3000  2981
2146
2147      P101 == 1.0
2148      1073  1074  1061  1056  1056
2149
2150      P101 == 2.0
2151      118   3133  3143  3125  3108

```

Listing 4.4: Example of preliminary system analysis

2152 The information obtained here can be used, as mentioned above, to
 2153 refine the enrichment of the dataset by setting directives in the [DATASET]
 2154 section of the *config.ini* file, should this be empty or only partially set, or to
 2155 make the first conjectures about the system, as we have just seen.

2156 The *system_info.py* file can also run in standalone mode if needed:
 2157 it takes as command-line arguments the dataset to be analyzed, a list of
 2158 actuators, and a list of sensors. For analysis related to state changes, the
 2159 dataset must mandatorily be of the enriched type.

2160 4.2.4 Phase 2: Graphs and Statistical Analysis

2161 The introduction of the new *graph analysis* is motivavted by from the re-
 2162 quirement to provide users with a comprehensive overview of the (sub)system
 2163 derived from the preceding pre-processing phase. The objective is to facil-
 2164 itate the identification of register types, enhance the understanding of re-
 2165 lationships, and effectively grasp the dynamics among registers controlled
 2166 by one or more PLCs. This analysis component serves to validate initial
 2167 conjectures made during the preliminary analysis described in the previ-
 2168 ous section or generate new insights with the aid of visual chart represen-
 2169 tations. It enables users to gain a deeper understanding of the system by
 2170 leveraging the support of visual charts.

2171 In Ceccato et al.'s framework, as mentioned in Section 3.2.7, it was only
 2172 possible to view the chart of one register at a time. While this allowed for
 2173 the identification or hypothesis of the register type, it made it challenging
 2174 to establish relationships with other components of the system and derive
 2175 conjectures about their behavior. To address this limitation, there was a
 2176 need for a new tool that could provide more information in a more acces-
 2177 sible manner.

2178 Initially, we considered adopting an approach similar to Figure 3.5,
 2179 where all the graphs are displayed within a single plot. However, we soon
 2180 realized that this solution was not feasible and could not be adopted. Fig-
 2181 ure 4.7 helps to understand why.

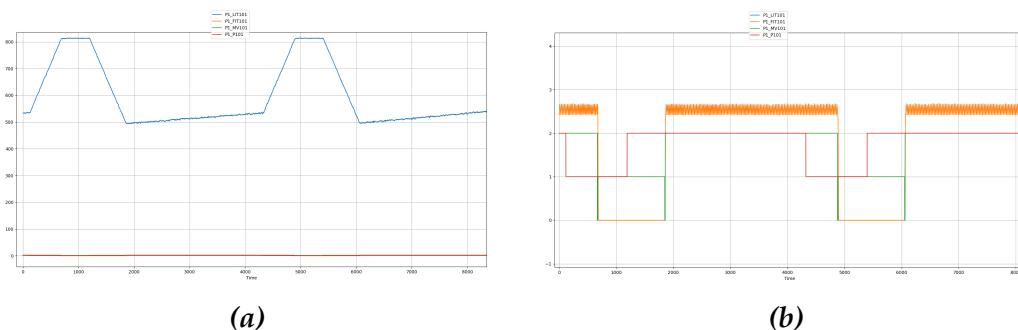
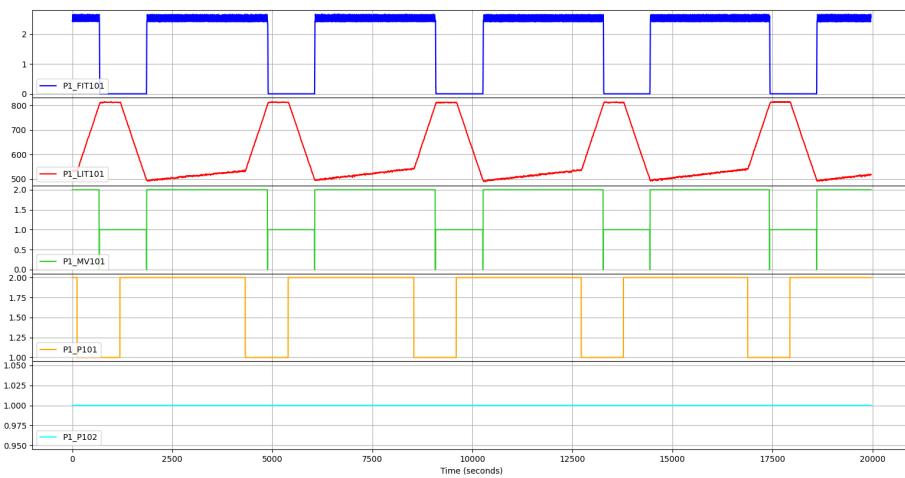


Figure 4.7: Plotting registers on the same y-axis

Figure 4.7a highlights the main issue with this approach, which is the use of the same y-axis for all the charts, representing the values of individual registers. When the range between register values is wide, it can result in some charts appearing as a single flat line or becoming indistinguishable, making them difficult to read. Additionally, as shown in Figure 4.7b, when registers have similar values, the graphs can become confusing and harder to interpret.

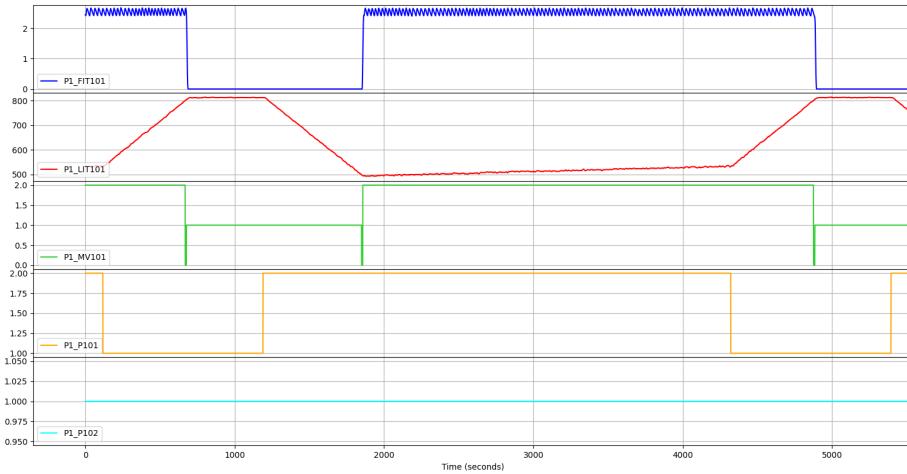
The solution to this issue is simple and effective: the use of **subplots**. Basically, each register corresponds to a subplot of the graph that shares the time axis (the x-axis) with the other subplots, but keeps the y-axis of the values of each register independent. This maintains the readability and comprehensibility of the charts, while simultaneously being able to immediately grasp the relationships between them. In addition, by sharing the time axis, it is possible to zoom in on a particular area of one of the charts and automatically the other ones will be zoomed in as well, thus not losing any information and no connection between registers.

Figure 4.8 illustrates more clearly what has just been explained.



(a) Example of plotting charts of a PLC registers using subplots

Figure 4.8: Example of the new graph analysis



(b) Zooming on a particular zone of the charts

Figure 4.8: Example of the new graph analysis (cont.)

2199 As the reader has probably already noticed, the majority of the graphs
 2200 presented in the previous sections and chapters were generated using the
 2201 new graph analysis script, specifically the `runChartsSubPlots.py` script.
 2202 This Python script is located in the `$(project-dir)/statistical-graphs`
 2203 directory and utilizes the `matplotlib` libraries to generate the graph plots,
 2204 similar to the Ceccato et al.'s tool.

2205 The script accepts the following command-line parameters:

- 2206 • **-f or --filename:** specifies the CVS format dataset to read data from.
 2207 The dataset must be within the directory containing the enriched
 2208 datasets for the invariant analysis phase. If subsequent parameters
 2209 are not specified, the script will display all registers in the dataset,
 2210 excluding any additional attributes;
- 2211 • **-r or --registers:** specifies one or more specific registers to be dis-
 2212 played;
- 2213 • **-a or --addregisters:** adds one or more registers to the default visual-
 2214 ization. This option is useful in case an additional attribute such as
 2215 slope is to be analyzed;

- 2216 • **-e or --excluderegisters:** excludes one or more specific registers from
2217 the default visualization. This option is useful to avoid displaying
2218 hardcoded setpoints or spare registers.

2219 This script, like the previous ones, is designed to provide the maximum
2220 flexibility and ease of use for the user, combined with greater power and
2221 effectiveness in deriving useful information about the analyzed system.

2222 **Statistical Analysis** After careful consideration, we made the decision
2223 not to include the statistical analysis aspect of the Ceccato et al.'s tool in
2224 the framework. We found that there was no practical use for it. Instead, we
2225 integrated the relevant statistical information into the preliminary analysis
2226 conducted after the pre-processing phase. Additionally, we deemed the
2227 histogram to have limited utility and considered it outdated in comparison
2228 to the new graph analysis approach we implemented.

2229 Although we decided not to include the statistical analysis aspect in
2230 the framework, the Python script `histPlots_Stats.py` from the original
2231 tool remains in the directory. This script is essentially unchanged from the
2232 version developed by Ceccato et al. and can be used in the future if the
2233 need arises.

2234 4.2.5 Phase 3: Invariant Inference and Analysis

2235 The phase of invariant inference and analysis has undergone a redesign
2236 and improvement to offer the user a more comprehensive and easier ap-
2237 proach to identify invariants. This has been achieved through the applica-
2238 tion of new criteria to analyze and reorganize the Daikon analysis results.
2239 The outcome of this is a more compact presentation of information that
2240 highlights the possible relationships among invariants.

2241 The new design not only enables the identification of undiscovered as-
2242 pects of the system behavior but also confirms the hypotheses made dur-
2243 ing the earlier stages of analysis. This step is now semi-automated, unlike
2244 before, and allows for the analysis of invariants on individual actuator
2245 states and their combinations.

2246 **4.2.5.1 Revised Daikon Output**

2247 To streamline the process of identifying invariants quickly and effi-
2248 ciently, it is necessary to revise the output generated by standard Daikon
2249 analysis. The goal is to create a more compact and readable format for the
2250 output.

2251 The current Daikon results basically consist of three sections, referred
2252 as *program point sections* [68]:

- 2253 1. the first section containing **general invariants**, i.e., valid regardless
2254 of whether a condition is specified for the analysis;
- 2255 2. the second section containing **conditional invariants**, obtained by
2256 specifying conditions for the analysis in the *.spininfo* file.
- 2257 3. a third section containing the invariants that are obtained from the
2258 **negation of the condition** potentially specified in the *.spininfo* file.

2259 In each section only a single invariant per row is shown, without relat-
2260 ing it in any way to the others: this makes it difficult to identify significant
2261 invariants and any invariant chain that might provide much more infor-
2262 mation about the behavior of the system than the single invariant.

2263 A brief example of the structure and format of this output related to
2264 PLC1 of the iTrust SWaT system is shown in Listing 4.5, where a condition
2265 was specified on the measurement LIT101 and on actuator MV101:

```
2266    aprogram.point:::POINT
2267     P102 == prev_P102
2268     FIT101 >= 0.0
2269     MV101 one of { 0.0, 1.0, 2.0 }
2270     P101 one of { 1.0, 2.0 }
2271     P102 == 1.0
2272     max_LIT101 == 816.0
2273     min_LIT101 == 489.0
2274     slope_LIT101 one of { -1.0, 0.0, 1.0 }
2275     [...]
2276     LIT101 > MV101
```

```

2277     LIT101 > P101
2278     LIT101 > P102
2279     LIT101 < max_LIT101
2280     LIT101 > min_LIT101
2281     [...]
2282     MV101 < min_LIT101
2283     MV101 < trend_LIT101
2284     P101 >= P102
2285     P101 < max_LIT101
2286     [...]
2287     =====
2288    aprogram.point:::POINT;condition="MV101 == 2.0 && LIT101 <
2289     ↪ max_LIT101 - 16 && LIT101 > min_LIT101 + 15"
2290     MV101 == prev_MV101
2291     P102 == slope_LIT101
2292     MV101 == 2.0
2293     FIT101 > MV101
2294     FIT101 > P101
2295     FIT101 > P102
2296     FIT101 > prev_P101
2297     MV101 >= P101
2298     MV101 >= prev_P101
2299     P101 <= prev_P101
2300     =====
2301    aprogram.point:::POINT;condition="not(MV101 == 2.0 &&
2302     ↪ LIT101 < max_LIT101 - 16 && LIT101 > min_LIT101 + 15)
2303     ↪ "
2304     P101 >= prev_P101
2305     Exiting Daikon.

```

Listing 4.5: Standard Daikon output for PLC1 of the iTrust SWaT system

2306 In the presented framework, the output is simplified to **two sections**:
2307 the *general section* and the section related to the *user-specified condition*. The
2308 section related to the negated condition is eliminated as it is not relevant in
2309 this context and could lead to potential misinterpretation. Moreover, the
2310 relationships between invariants will be emphasized by utilizing **transi-**
2311 **tive closures**: transitive closure of a relation R is another relation, typically
2312 denoted R^+ that adds to R all those elements that, while not necessarily

2313 related directly to each other, can be reached by a *chain* of elements related
2314 to each other. In other words, the transitive closure of R is the smallest (in
2315 set theory sense) transitive relation R such that $R \subset R^+$ [69].

To implement the transitive closure of the invariants generated by Daikon, we initially categorized the invariants in each section by type based on their mathematical relation ($==$, $>$, $<$, $>=$, $<=$, $!=$), excluding any invariant related to additional attributes except for the slope. For each type of invariant, we constructed a **graph** using the NetworkX library, where registers were represented as nodes, and arcs were created to connect registers that shared a common endpoint in the invariant, applying the transitive property. To reconstruct the individual invariant chains, we employed a straightforward approach known as *Depth-first Search* (DFS) on each of the graphs. This method allowed us to traverse the graphs and obtain the desired outcome, identifying the invariant chains. Figure 4.9 shows some examples of these graphs:

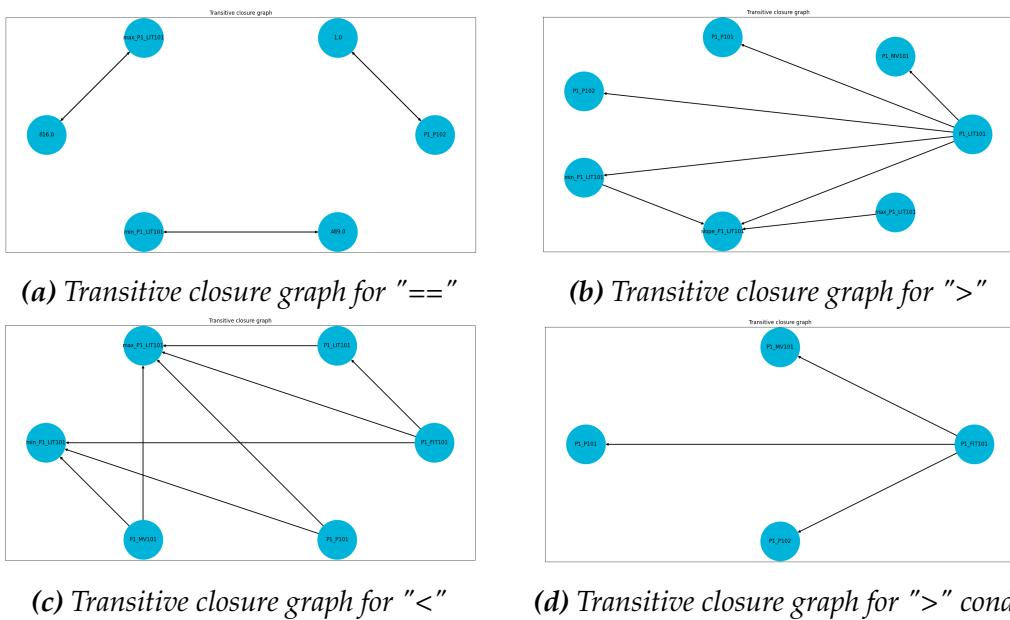


Figure 4.9: Example of transitive closure graphs for invariants in PLC1 of iTrust SWaT system

At the end of this process, still applied to PLC1 of the iTrust SWaT sys-

2329 tem and with the same analysis condition, we get the following complete
 2330 output:

```

2331 1 =====
2332 2 General
2333 3 =====
2334 4 MV101 one of { 0.0, 1.0, 2.0 }
2335 5 P101 one of { 1.0, 2.0 }
2336 6 slope_LIT101 one of { -1.0, 0.0, 1.0 }
2337 7 FIT101 != P101, P102
2338 8 P102 == 1.0
2339 9 max_LIT101 == 816.0
2340 10 min_LIT101 == 489.0
2341 11 LIT101 > MV101
2342 12 LIT101 > P101
2343 13 LIT101 > P102
2344 14 LIT101 > min_LIT101 > slope_LIT101
2345 15 FIT101 >= 0.0
2346 16 P101 >= P102 >= slope_LIT101
2347 17
2348 18 =====
2349 19 MV101 == 2.0 && LIT101 < max_LIT101 - 16 && LIT101 >
2350 20 ↢ min_LIT101 + 15
2351 21 =====
2352 22 slope_LIT101 == P102
2353 23 MV101 == 2.0
2354 24 FIT101 > MV101
2355 25 FIT101 > P101
2356 26 FIT101 > P102
2357 27 MV101 >= P101
  
```

***Listing 4.6:** Revised Daikon output with transitive closures for PLC1 of the iTrust SWaT system*

2358 Transitive closures can be appreciated in lines 7, 14 and 16 of Listing 4.6.
 2359 In general, the output has been reduced in the number of effective rows
 2360 (19 versus the 61 in the output of Listing 4.5, making it certainly better to
 2361 read and identify significant invariants) and the invariant chains make it
 2362 more immediate to grasp the relationships between registers.

2363 **4.2.5.2 Types of Analysis**

2364 In contrast to Ceccato et al.'s solution, which involves manual individual analyses, our proposal is to introduce **two types of semi-automated analysis**.

2367 The first type focuses on analyzing **all states for each individual actuator**. This automated analysis aims to provide comprehensive insights into the behavior of each actuator in relation to a specific measurement selected by the user.

2371 The second type of analysis considers the current system configuration and examines the **actual states of the actuators**. This analysis takes into account the interplay and combined effects of multiple actuators on the selected measurement. By incorporating the real-time states of the actuators, this semi-automated analysis offers a more accurate assessment of the system behavior.

2377 These two types of analysis will be handled by the Python script `daikonAnalysis.py`, contained in the default directory `$(project_dir)/daikon`.
2378
2379 The script accepts the following command-line arguments:

- 2380 • **-f or --filename**: specifies the enriched dataset, in CSV format, from
2381 which to read the data. The dataset must be located within the directory
2382 containing the enriched datasets specified in the `config.ini` file
2383 (by default `$(project_dir)/daikon/Daikon_Invariants`);
- 2384 • **-s or --simpleanalysis**: performs the analysis on the states of individual
2385 actuators;
- 2386 • **-c or --customanalysis**: performs the analysis on combinations of actual
2387 states of the actuators;
- 2388 • **-u or --uppermargin**: defines a percentage margin on the maximum
2389 value of the measurement;
- 2390 • **-l or --lowermargin**: defines a percentage margin on the minimum
2391 value of the measurement;

2392 The selection of one or both types of analysis is possible. Additionally,
 2393 the last two parameters can be used to set a condition on the value of the
 2394 measurement. This condition is designed to bypass the transient periods
 2395 that occur during actuator state changes and the actual trend changes at
 2396 the maximum and minimum values of the measurement.

2397 This approach proves particularly beneficial for the first type of analy-
 2398 sis, as it enables more accurate data on the trends of the measurement. By
 2399 excluding the transient periods, the analysis can focus on the stable and
 2400 meaningful trends, providing improved insights into the behavior of the
 2401 system.

2402 **Analysis on single actuator states** Analysis on the states of individual
 2403 actuators is the simplest: after the user is prompted to input the measure-
 2404 ment, chosen from a list of likely available measurements, the script rec-
 2405 ognizes the likely actuators and the relative states of each, using the same
 2406 mixed Daikon/Pandas technique adopted in the preliminary analysis dur-
 2407 ing the pre-processing phase.

2408 For each actuator and each state it assumes, a single Daikon analysis is
 2409 performed, eventually placing the condition on the maximum and mini-
 2410 mum level of the measurement.

2411 The result of these analyses are saved in the form of text files in a di-
 2412 rectory having the name corresponding to the analyzed actuator and con-
 2413 tained in the default parent directory

2414 \$(project_dir)/daikon/Daikon_Invariants/results: each file generated
 2415 by the analysis is identified by the name of the actuator, the state and the
 2416 condition, if any, on the measurement (see Figure 4.10).

2417 Listing 4.6 provides an illustrative example of the analysis results. In
 2418 this case, we focus on the actuator MV101 in state 2.

2419 The condition-generated invariants provide the most interesting in-
 2420 sights. For example, at line 21, we observe that when MV101 is set to 2,
 2421 the slope of LIT101 is 1, indicating an increasing trend. This leads us to
 2422 speculate that MV101 represents the actuator's ON state and is responsible
 2423 for filling the tank (LIT101).

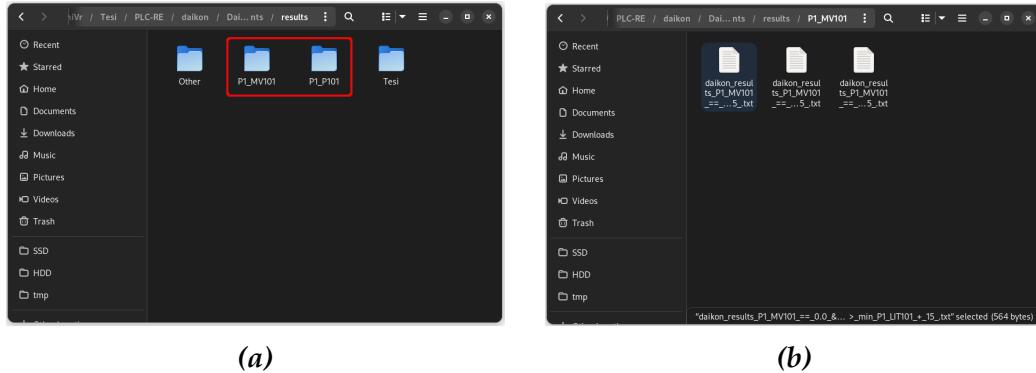


Figure 4.10: Directory (a) and outcome files (b) for the single actuator states analysis

2424 **Analysis of the Current System Configuration** The analysis of the cur-
 2425 rent system configuration based on the actual states of the actuators is
 2426 more complex, but at the same time offers more interesting outcomes as
 2427 it provides better evidence of actuator behavior in relation to the selected
 2428 measurement.

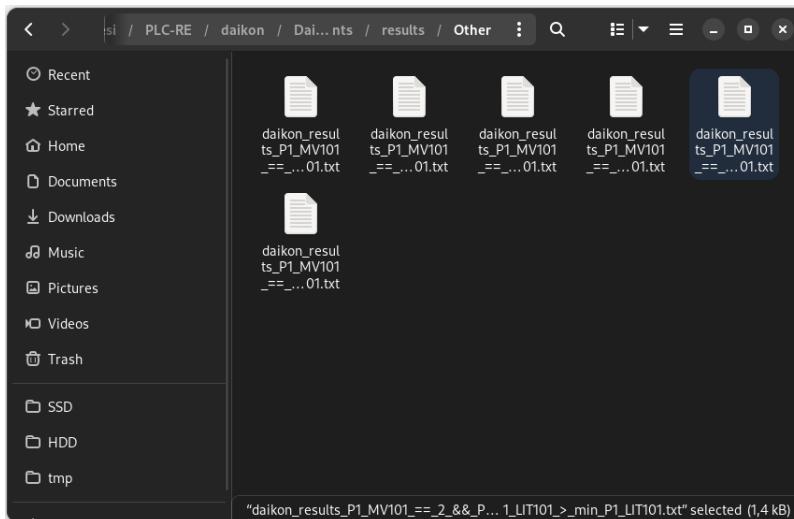


Figure 4.11: Daikon outcome files for system configuration analysis. Each file represents a single system state

2429 For this analysis, the script automatically identifies the configurations of
 2430 the actuators that represent different system states (e.g., MV101 == 2, P101
 2431 == 1). Daikon analysis is then performed for each of these configurations.

2432 The user is prompted to select the measurement attribute and, if desired,
2433 specific actuators for studying their configurations. If no actuators are se-
2434 lected, all previously detected actuators will be considered.

2435 The analysis results are saved in text format within a designated di-
2436 rectory, located alongside the previous analysis outputs. The filenames of
2437 these result files follow a specific naming convention based on the analysis
2438 rule applied (see Figure 4.11).

2439 An example of the obtained outcomes can be seen in Listing 4.7:

```

2440   1 =====
2441   2 General
2442   3 =====
2443   4 MV101 <= P101 ==> FIT101 >= 0.0
2444   5 MV101 <= P101 ==> MV101 one of { 0.0, 1.0, 2.0 }
2445   6 MV101 <= P101 ==> P101 one of { 1.0, 2.0 }
2446   7 MV101 <= P101 ==> slope_LIT101 one of { -1.0, 0.0, 1.0 }
2447   8 MV101 > P101 ==> FIT101 > MV101
2448   9 MV101 > P101 ==> FIT101 > P101
2449  10 MV101 > P101 ==> FIT101 > P102
2450  11 MV101 > P101 ==> FIT101 > slope_LIT101
2451  12 MV101 > P101 ==> MV101 == 2.0
2452  13 MV101 > P101 ==> MV101 > P102
2453  14 MV101 > P101 ==> MV101 > slope_LIT101
2454  15 MV101 > P101 ==> P101 == 1.0
2455  16 MV101 > P101 ==> P101 == P102
2456  17 MV101 > P101 ==> P101 == slope_LIT101
2457  18 MV101 > P101 ==> slope_LIT101 == 1.0
2458  19 [...]
2459  20
2460  21 =====
2461  22 MV101 == 2 && P101 == 1 && LIT101 < max_LIT101 && LIT101 >
2462    ↢ min_LIT101
2463  23 =====
2464  24 slope_LIT101 == P102 == P101 == 1.0
2465  25 MV101 == 2.0
2466  26 FIT101 > MV101

```

2467 27 FIT101 > P101

Listing 4.7: Daikon outcomes for the system configuration MV101 == 2, P101 == 1 on LIT101

2468 In contrast to Listing 4.6, the current analysis reveals the presence of impli-
2469 cations in the general invariants section, which were previously absent (re-
2470 maining generic invariants omitted for brevity). These implications offer
2471 valuable insights, as demonstrated in this case. For instance, the invariant
2472 stated in line 18 informs us that if the value of MV101 is greater than P101,
2473 then the slope's value is 1, indicating an increasing tank level. This finding
2474 further corroborates our previous understanding that the state MV101 ==
2475 2 represents the ON state for the actuator responsible for filling the tank,
2476 namely LIT101.

2477 **Refining the Analysis** In certain situations, the outcomes provided by
2478 the semi-automated analyses may not meet the user's expectations. For
2479 instance, the clarity of the slope value may be insufficient, or the user may
2480 wish to delve deeper into a specific aspect of the system to uncover ad-
2481 dditional invariants that were not previously identified. In such cases, the
2482 user has the option to conduct a more targeted and specific invariant anal-
2483 ysis using the Python script `runDaikon.py`, which enables precise investi-
2484 gations of the system.

2485 The script, located in the default directory `$(project_dir)/daikon`, can
2486 be executed with three command-line parameters:

- 2487 • **-f or --filename:** specifies the CSV format *enriched* dataset to read
2488 data from. Even in this case, the dataset must be located within the
2489 directory containing the enriched datasets;
- 2490 • **-c or --condition:** specifies the condition for the analysis, which will
2491 be automatically overwritten in the `.spinfo` file. It is possible to spec-
2492 ify more than one condition, but it is strongly recommended to use
2493 the logical operator `&&` to avoid undesired behaviors in Daikon out-
2494 comes;

- 2495 • **-r or --register:** specifies the directory where to save the text file with
2496 the outcomes.

2497 During the execution of this analysis, a single output file is generated,
2498 containing the discovered invariants. The user can specify the directory
2499 where this file should be saved using the `-r` command-line option. By de-
2500 fault, the file is stored in the directory
2501 `$(project_dir)/daikon/Daikon_Invariants/results`. The output file serves
2502 as a record of the identified invariants and can be examined at a later time.

2503 In conclusion, the integration of these two analysis types, along with
2504 the ability to conduct more refined analysis in the future and the enhanced
2505 output format of Daikon, significantly enhances the completeness, clarity,
2506 and effectiveness of this stage compared to the Ceccato et al.'s framework.

2507 4.2.6 Phase 4: Business Process Analysis

2508 We have made significant revisions to the Business Process Analysis
2509 we are presenting. Instead of relying on the previous Java solution and
2510 proprietary process mining software Disco, we have adopted a **new inte-**
2511 **grated solution** created in Python from scratch. The new solution utilizes
2512 the Graphviz libraries to generate the corresponding activity diagram.

2513 In this updated Business Process, greater emphasis is placed on pro-
2514 cess mining related to the physical system. Our goal is to extract as much
2515 information as possible from the dataset, enabling us to promptly visu-
2516 alize the system's behavior and its various states. This approach allows
2517 us to validate the conjectures and hypotheses formulated in the previ-
2518 ous phases, and potentially uncover hidden patterns that were previously
2519 undisclosed.

2520 On the other hand, the aspect related to network communications was
2521 reconsidered to enable operation with multiple protocols, expanding be-
2522 yond the limitations of Modbus.

2523
2524 Now, let's examine the key aspects of this new phase in greater detail.

2525 4.2.6.1 Process Mining of the Physical Process

2526 The mining of the physical process is performed by the Python script
2527 called `processMining.py`, located in the default directory `$(project_dir)/process-mining`.
2528 This script accepts the following parameters from the command line:

- 2529 • **-f or --filename:** specifies the CSV format *timestamped* dataset to be
2530 mined from. The dataset is obtained from the pre-processing stage
2531 and is located in the default directory `$(project_dir)/process-mining/data`;
- 2532 • **-a or --actuators:** specifies one ore more actuators whose combina-
2533 tions of states are to be analyzed. If this parameter is not provided,
2534 all actuators in the subsystem will be considered;
- 2535 • **-s or --sensors:** specifies one or more measurements for which the
2536 trend will be calculated based on actuator state changes. If this pa-
2537 rameter is omitted, all available measurements will be considered;
- 2538 • **-t or --tolerance:** specifies the tolerance to be taken into account dur-
2539 ing the trend calculation;
- 2540 • **-g or --graph:** shows the resulting activity diagram.

2541 The script processes the dataset in a sequential manner, analyzing each
2542 actuator state change. It calculates the duration in seconds of the system
2543 state, as well as the trend and slope of the specified measurement(s). Addi-
2544 tionally, the script stores the next system state and the measurement values
2545 corresponding to the state change points, allowing for the identification of
2546 relative setpoints. These results are gradually collected and stored in a dic-
2547 tionary, where the keys represent the system states based on the actuator
2548 configurations, and the values represent the measured values mentioned
2549 above. The dictionary is then saved as a JSON file, which can be accessed
2550 by the user in the `$(project_dir)/process-mining/data` directory. The
2551 name of the file can be specified in the configuration file `config.ini`.

2552 An example of the JSON file obtained in this step is shown in Listing
2553 4.8: the JSON file showcases the structure of the data obtained during the
2554 process.

```
2555 {
2556     "MV101 == 2, P101 == 2": {
2557         "start_value_LIT101": [
2558             534.9366,
2559             495.4483,
2560             497.9998,
2561             495.9586,
2562             495.8016
2563         ],
2564         "end_value_LIT101": [
2565             536.0356,
2566             532.7384,
2567             541.7273,
2568             534.4656,
2569             540.8245
2570         ],
2571         "slope_LIT101": [
2572             0,
2573             0.015,
2574             0.018,
2575             0.016,
2576             0.018
2577         ],
2578         "trend_LIT101": [
2579             "STBL",
2580             "ASC",
2581             "ASC",
2582             "ASC",
2583             "ASC"
2584         ],
2585         "time": [
2586             119,
2587             2464,
2588             2477,
2589             2452,
2590             2440
```

```

2591     ],
2592     "next_state": [
2593       "MV101 == 2, P101 == 1",
2594       "MV101 == 2, P101 == 1",
2595       "MV101 == 2, P101 == 1",
2596       "MV101 == 2, P101 == 1",
2597       "MV101 == 2, P101 == 1"
2598     ]
2599   },
2600   "MV101 == 2, P101 == 1": {
2601     ...
2602   }
2603   "MV101 == 0, P101 == 1": {
2604     ...
2605   }
2606   ...
2607 }
2608 }
```

Listing 4.8: Example of the data contained in the produced JSON file

The collected data is now utilized to generate the **system activity diagram**. This diagram represents an *oriented* graph where the nodes correspond to the **system states** determined by the actuator configurations. In addition to the state information, the nodes also display the **trend** (ascending, descending, or stable) and the slope of the reference measurement.

The edges in the diagram depict the values of the measurements at the time of the state change. These measurement values represent the **absolute setpoints** for each measurement. Setpoints are calculated on the average of the values measured for that specific state. Additionally, the diagram incorporates on the edges the average duration of each system state.

Each data point on the edges is accompanied by a *standard deviation value*. This value provides information about the occurrence of a specific state within the system's cycle, indicating whether it appears multiple times with varying values and time durations. A low standard deviation suggests that the state is likely to occur only once within each cycle.

2625 Conversely, a high standard deviation indicates that the state may occur
 2626 multiple times within the cycle.

2627 An example of the activity diagram generated by the processMining.py
 2628 script is presented in Figure 4.12. This diagram illustrates the system's
 2629 behavior by depicting transitions between different states. Nodes in the
 2630 diagram represent specific system states, while arrows or edges indicate
 2631 the flow between these states.

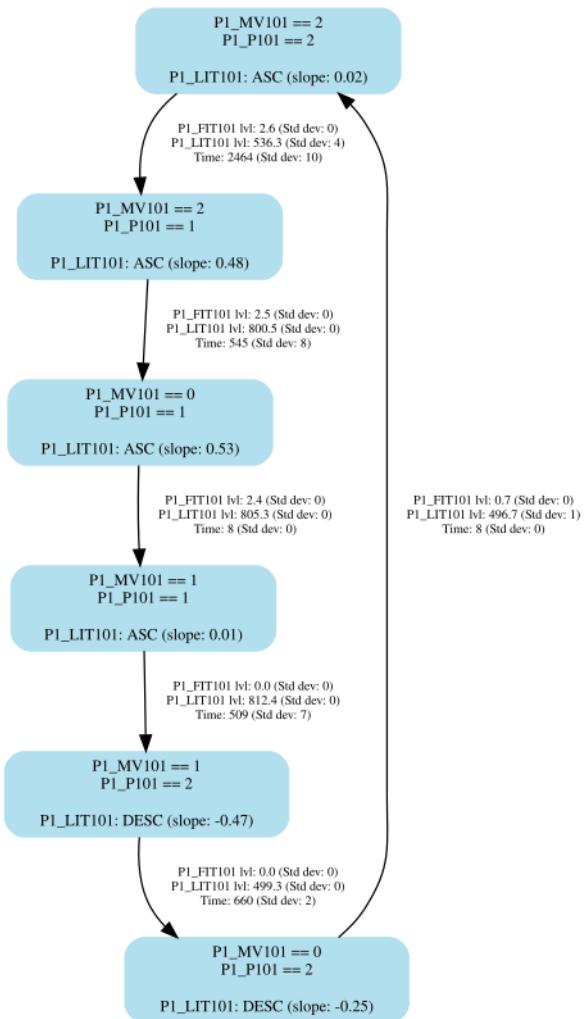


Figure 4.12: Activity diagram for PLC1 of the iTrust SWaT system

2632 The diagram reveals important aspects of the system, such as its **periodic-**

2633 ity and the emphasis on the **temporal sequence of actuator states**. These
2634 insights might not have been clearly evident in earlier stages of the anal-
2635 ysis. The diagram provides a visual representation that allows us to ap-
2636 preciate the recurring patterns and understand how the system evolves
2637 over time. By observing the transitions and relationships between differ-
2638 ent states, we gain a better understanding of the dynamics and behavior
2639 of the system.

2640 It is important to acknowledge that despite considering the tolerance
2641 set for trend and slope calculation, the accuracy of the related data may
2642 vary. For instance, there could be instances where multiple trends exist
2643 within the same node, or a trend may be stable instead of increasing or de-
2644 creasing. These discrepancies arise due to **data perturbations**, which were
2645 discussed in Section 4.2.3.2. Additionally, the noise reduction techniques
2646 seen in the same Section were not employed in this analysis. Therefore, it
2647 becomes the responsibility of the user to correctly interpret the outcomes
2648 of this step, taking into account the information presented in the process
2649 mining diagram and the findings from previous analyses. By consider-
2650 ing these factors together, the user can make informed interpretations and
2651 draw accurate conclusions from the results.

2652 **4.2.6.2 Network Data**

2653 The incorporation of network data into Business Process Analysis has
2654 been reconsidered to shift away from a *single-protocol* solution based on
2655 Modbus. Instead, the focus is on adopting a solution that can handle **mul-**
2656 **tiple protocols**, even at the same time.

2657 The main concept was to develop a new Python script that would ex-
2658 tract and process data from PCAP files obtained through network traffic
2659 sniffing. The intention was to export this data to a CSV file, which would
2660 then be used as input for the process mining script, `processMining.py`,
2661 and integrated with the physical system data to derive commands to the
2662 actuators.

2663 The analysis of the network data revealed that different protocols ex-
 2664 hibit distinct behaviors when commanding changes in actuator states. Con-
 2665 sequently, the previous approach proposed by Ceccato et al. becomes **im-**
 2666 **practical** when attempting to detect system state changes through network
 2667 commands sent to the actuators.

2668 However, considering that system state changes have already been
 2669 identified through the process mining of the physical process, and with
 2670 the corresponding event timestamps available, we propose an alternative
 2671 approach to Ceccato et al.'s method. Instead of seeking the correspon-
 2672 dence between network events and physical process events at the same
 2673 instant, we suggest reversing the perspective. By focusing on the corre-
 2674 spondence between a given event occurring in the physical process at a
 2675 specific moment and the corresponding event in the network data at the
 2676 same moment, it should be possible to achieve a similar, if not superior,
 2677 outcome compared to the previous solution.

2678

4.2.7 Summary

2679 In Table 4.1, we provide a *phase-by-phase* summary of the enhance-
 2680 ments achieved by the proposed framework compared to the limitations
 2681 observed in Ceccato et al., as presented in Table 3.1.

Phase	Enhancements
General	- Independence from the analyzed system achieved
Improvements	<ul style="list-style-type: none"> through a general <i>config.ini</i> configuration file. - Command line functionality improved for enhanced usability. - Utilization of a single programming language throughout the implementation.

Network Analysis	<ul style="list-style-type: none">- Introduction of a novel phase, not found in Ceccato et al.- Reconstruction of the network topology and communication links between PLCs, as well as between PLCs and other devices such as HMI and Historian Server.- Identification of industrial protocols used in the system.- Provision of graphical and textual output representing the network topology and communication details.
Pre-processing	<ul style="list-style-type: none">- Capability to choose a subsystem by specifying individual PLCs and the desired time range for analysis.- Enhanced and automated dataset enrichment that is optimized and customizable, enabling the allocation of additional fields to specific registers or groups of registers.- Mitigation of noise caused by data perturbations through the application of Seasonal and Trend decomposition using Loess (STL).
Preliminary Analysis	<ul style="list-style-type: none">- Introduction of a novel feature, not included in Ceccato et al.- Automatic identification of potential actuators, measurements and hardcoded setpoints/spare actuators.- Analysis of state transitions of individual actuators in correlation with a specific measurement.- Time duration analysis of actuator states.
Graphical / Statistical Analysis	<ul style="list-style-type: none">- Creating visual representations of PLC registers through the use of subplots instead of single plots.- Providing a comprehensive view of the system.- It is possible to zoom in on a particular area of one of the charts and automatically the other ones will be zoomed in as well, thus not losing any information and no connection between registers.- Enabling the selection of specific registers for analysis.

Invariant Analysis	<ul style="list-style-type: none"> - Revised Daikon output through the utilization of transitive closures, resulting in more concise and readable results. - Enhancing the identification of invariants by making the output easier to interpret. - Semi-automated generation of invariants by leveraging identified actuators.
Business Process Analysis	<ul style="list-style-type: none"> - Complete overhaul of the phase, developed from scratch without reliance on external tools. - Heightened focus on the physical process aspect of the system. - Recognition of actual states of the system, considering their trends in relation to a specific measure. - Detection of changes in slopes within the same trend. - Determination of absolute setpoints with respect to measurements. - Introduction of a new system activity diagram to depict system behavior.

Table 4.1: Summary table of proposed framework enhancements

Case study: the iTrust SWaT System

2682 HAVING introduced the innovative framework and highlighted its po-
2683 tential in the preceding chapter, we now turn our attention to the
2684 case study where we will apply this framework. As previously mentioned
2685 in Chapter 4 and demonstrated through various examples in the same
2686 chapter, our focus will be on the **iTrust SWaT system** [15], developed by
2687 the iTrust – Center for Research in Cyber Security of University of Singa-
2688 pore for Technology and Design [37]. The acronym SWaT represents *Secure*
2689 *Water Treatment*.

2690 The iTrust SWaT system is a testbed that replicates on a small scale
2691 a real water treatment plant arises to support research in the area of cy-
2692 ber security of industrial control systems and has been operational since
2693 March 2015: it is still being used by students at the University of Singapore
2694 for educational and training purposes and is available to organizations to
2695 train their operators on cyber physical incidents.

2696 5.1 Architecture

2697 In contrast to the virtualized testbed discussed in Section 3.2.1 by Cec-
2698 cato et al., the iTrust SWaT system is composed entirely of physical hard-
2699 ware components. It encompasses various elements, starting from field

2700 devices and extending to PLCs, HMI, SCADA workstations, and the SCADA
2701 server (also referred to as the *historian*). The historian is responsible for
2702 recording data from field devices for further analysis. In the upcoming
2703 sections, we will delve deeper into the architecture of the physical process
2704 and the communication network.

2705 **5.1.1 Physical Process**

2706 The physical process of the SWaT consists of six stages, denoted P1
2707 through P6 [70][71]. These stages are:

- 2708 P1. **taking in raw water:** feeds unfiltered water into the system
- 2709 P2. **chemical dosing:** adds chemicals to water useful for initial pretreat-
2710 ment;
- 2711 P3. **Ultra Filtration (UF) system:** the water is filtered through a semi-
2712 permeable membrane (ultrafiltration membrane) using the liquid pres-
2713 sure, effectively capturing impurities and suspended solids, as well
2714 as removing bacteria, viruses, and other pathogens present in the
2715 water;
- 2716 P4. **dechlorination:** removes residual chlorine from disinfected water
2717 using ultraviolet lamps;
- 2718 P5. **Reverse Osmosis (RO):** performs further filtration of the water;
- 2719 P6. **backwash process:** cleans the membranes in UF using the water pro-
2720 duced by RO.

2721 Figure 5.1 [15] shows a graphical representation of the architecture and the
2722 six stages of the SWaT system.

2723 The SWaT system incorporates an array of sensors that play a crucial
2724 role in monitoring the system's operations and ensuring their safety. These
2725 sensors are responsible for continuously collecting data and providing in-
2726 teresting insights into the functioning of the system [70]. These sensors
2727 are:

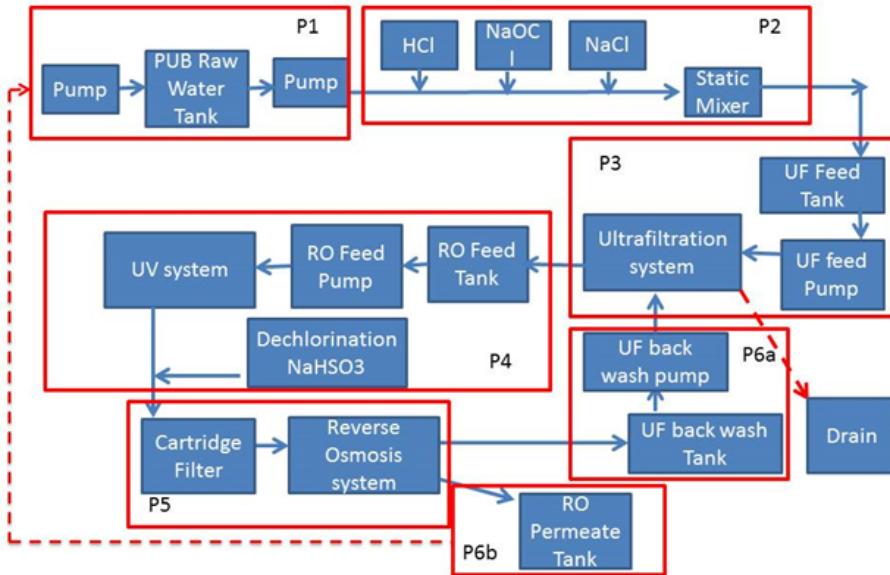


Figure 5.1: SWaT architecture

- Level Indication Transmitter (measured in mm);
- Flow Indication Transmitter (m³/hr);
- Analyser Indicator Transmitter;
 - Conductivity ($\mu\text{S}/\text{cm}$);
 - pH;
 - Oxidation Reduction Potential (mV);
- Differential Pressure Indicator Transmitter (kPa);
- Pressure Indicator Transmitter (kPa);

Sensors and actuators associated with each PLC are shown in Figure 5.2 [70]. Sensors and actuators are mapped into tags by the communication protocol used (see 5.1.2): a tag can be addressed via string descriptor defined by the system designer (e.g. MV101, to indicate motorized valve number 1 at stage 1) or by referring directly to the analog/digital pins of the PLC I/O unit [71].

Raw Water	Pre-Treatment	Ultra-Filtration	De-Chlorination	Reverse Osmosis	RO Product
P-101 Stopped	P-201 Stopped	P-301 Stopped	P-401 Stopped	P-501 Stopped	P-601 Stopped
P-102 Stopped	P-202 Stopped	P-302 Stopped	P-402 Stopped	P-502 Stopped	P-602 Stopped
MV-101 Closed	P-203 Stopped	MV-301 Closed	P-403 Stopped	MV-501 Closed	P-603 Stopped
LIT-101 520 mm	P-204 Stopped	MV-302 Closed	P-404 Stopped	MV-502 Closed	LS-601 Normal
FIT-101 0.00 m³/h LL	P-205 Stopped	MV-303 Closed	UV-401 Stopped	MV-503 Closed	LS-602 HIGH
FIT-201 0.00 m³/h LL	P-206 Stopped	MV-304 Closed	LS-401 Normal	MV-504 Closed	LS-603 LOW
	P-207 Stopped	PSH-301 Normal	LIT-401 1008 mm H	PSL-501 Normal	FIT-601 0.00 m³/h LL
	P-208 Stopped	DPSH-301 Normal	FIT-401 0.00 m³/h LL	PSH-501 Normal	
	MV-201 Closed	LIT-301 1012 mm H	AIT-401 0.17 ppm	AUT-501 6.89	
	LS-201 Normal	FIT-301 0.00 m³/h	AIT-402 275.70 mV	AUT-502 204.20 mV	
	LS-202 Normal	DPT-301 0.95 kPa LL		AUT-503 264.23 µS/cm H	
	LS-203 Normal			AUT-504 14.27 µS/cm H	
	AIT-201 142.18 µS/cm L			FIT-501 0.00 m³/h L	
	AIT-202 7.20 H			FIT-502 0.00 m³/h HH	
	AIT-203 293.59 mV L			FIT-503 0.00 m³/h HH	
				FIT-504 0.00 m³/h LL	
				PIT-501 2.64 kPa LL	
				PIT-502 0.00 kPa H	
				PIT-503 0.00 kPa	

Figure 5.2: Sensors and actuators associated with each PLC

2742 5.1.2 Control and Communication Network

2743 The SWaT system's network architecture follows the principles of lay-
 2744 ering and zoning, which enable segmentation and control of traffic within
 2745 the network.

2746

2747 Five layers are present starting from the highest to the lowest:

- 2748 • Layer 3.5 – Demilitarized Zone (DMZ);
- 2749 • Layer 3 – Operation Management (Historian);
- 2750 • Layer 2 – Supervisory Control (Touch Panel, Engineering Worksta-
2751 tion, HMI Control Clients);
- 2752 • Layer 1 – Plant Control Network (PLCs) (Star Network);
- 2753 • Layer 0 – Process (Actuator/Sensors and Input/output modules)
2754 (Ring Network).

2755 PLCs at Layer 1 communicate with their respective sensors and actu-
 2756 ators at Layer 0 through a conventional ring network topology based on

2757 EtherNet/IP, to ensure that the system can tolerate the loss of a single link
2758 without any adverse impact on data or control functionality.

2759 PLCs between the different process stages at Layer 1 communicate
2760 with each other through a star network topology using the CIP protocol
2761 on EtherNet/IP, previously discussed in Section 2.2.6.3.

2762 Regarding zoning, the SWaT system is divided into three zones, each
2763 containing one or more layers. These zones are, in descending order of
2764 security level:

- 2765 • **Plant Control Network, or Control System:** includes layers from 0
2766 to 2;
- 2767 • **DMZ:** includes Layer 3.5;
- 2768 • **Plant Network:** includes Layer 3;

2769 Figure 5.3 [15] provides a clearer visualization of the zoning and layer
2770 division within the network architecture of the SWaT system. This dia-
2771 gram highlights the distinct zones and their corresponding layers, offering
2772 a comprehensive overview of the system's network structure.

2773 A specific IP address is associated with each device: in Table 5.1 we
2774 report the addresses for the PLCs, historian, and Touch Panel in the *Plant*
2775 *Control Network* (PCN) zone:

IP Address	Device
192.168.1.10	PLC1
192.168.1.20	PLC2
192.168.1.30	PLC3
192.168.1.40	PLC4
192.168.1.50	PLC5
192.168.1.60	PLC6
192.168.1.100	Touch Panel (PCN)
192.168.1.200	Historian Server

192.168.1.201	Engineering Workstation
---------------	-------------------------

Table 5.1: Main IP addresses of the six PLCs and SCADA in the SWaT system

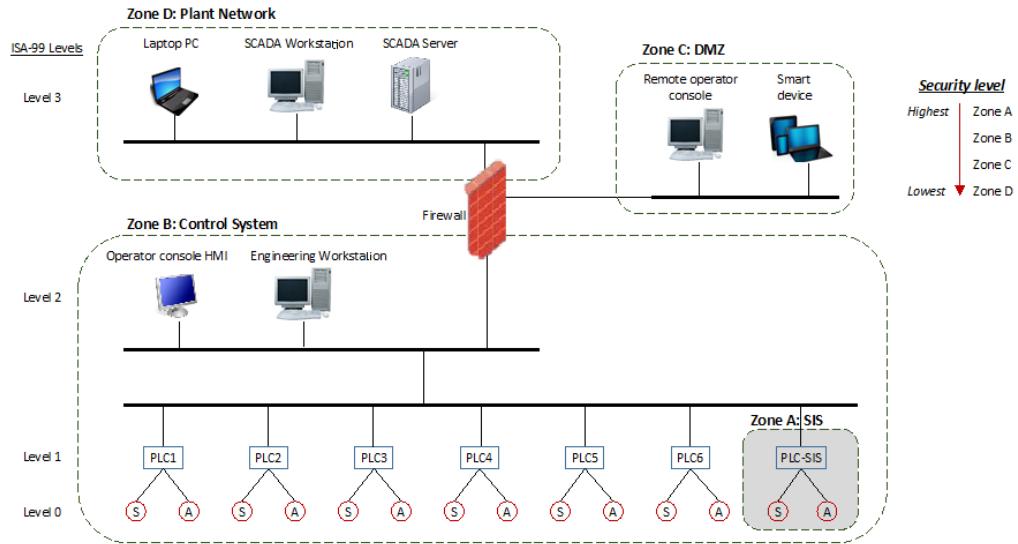


Figure 5.3: SWaT network architecture and zoning

2776 5.2 Datasets

2777 To facilitate the study and testing of technologies related to cyber secu-
 2778 rity in Industrial Control Systems and critical infrastructure, iTrust offers
 2779 researchers worldwide the opportunity to access a range of datasets [72].
 2780 These datasets consist of a collection of data obtained from the SWaT sys-
 2781 tem, encompassing information on both the physical processes and net-
 2782 work communications. The data is organized into different years, and
 2783 researchers can request access to these datasets for their analysis and ex-
 2784 perimentation purposes.

2785 **Physical Process Datasets** The datasets containing information about
 2786 the physical processes are provided in CSV format files. These files en-

2787 compass data collected during different time intervals, which can vary
2788 from a few hours to entire days. The granularity of the data is typically
2789 at a one-second interval, although there may be some exceptions. The col-
2790 lected data primarily consists of timestamped sensor measurements and
2791 actuator status values for each PLC, describing the physical properties of
2792 the testbed in operational mode.

2793 **Network Communications Datasets** Network communications are typi-
2794 cally available in the form of *Packet Capture* format (PCAP) files. These files
2795 contain captures of communication network traffic, allowing researchers
2796 to analyze and examine the network interactions. In some instances, CSV
2797 files are provided instead of PCAP files, featuring different characteristics
2798 for the collected data.

2799 5.2.1 Our Case Study: the 2015 Dataset

2800 The dataset selected as a case study to apply the framework discussed
2801 in the previous chapter is specifically the dataset from the year 2015 [73].
2802 The main reason for this choice is the unique characteristics found in the
2803 physical process dataset that are not present in datasets from subsequent
2804 years.

2805 **Physical Process Data** The data collection process lasted 11 consecutive
2806 days, 24 hours per day. During the first 7 days, the system operated nor-
2807 mally without any recorded attacks. However, attacks were observed dur-
2808 ing the remaining 4 days. The collected data reflects the impact of these
2809 attacks, leading to the creation of two separate CSV files: one containing
2810 the recorded data of SWaT during the system's regular operations, and
2811 the other containing data recorded during the days of the attacks. To en-
2812 sure accurate information about the system, the dataset pertaining to the
2813 normal operations, which spans seven days, was chosen for analysis.

2814 Data collection occurs at a frequency of one data point per second,
2815 with the assumption that significant attacks cannot occur within a shorter

2816 time frame. Additionally, the firmware of the PLCs remains unchanged
 2817 throughout the data collection period.

2818 At the beginning of data gathering the tanks are empty and the system
 2819 must be initialized in order to then reach full operation: it typically takes
 2820 around five hours for all tanks to be fully filled and for the system to sta-
 2821 bilize and reach the appropriate operational state (see Figure 5.4).

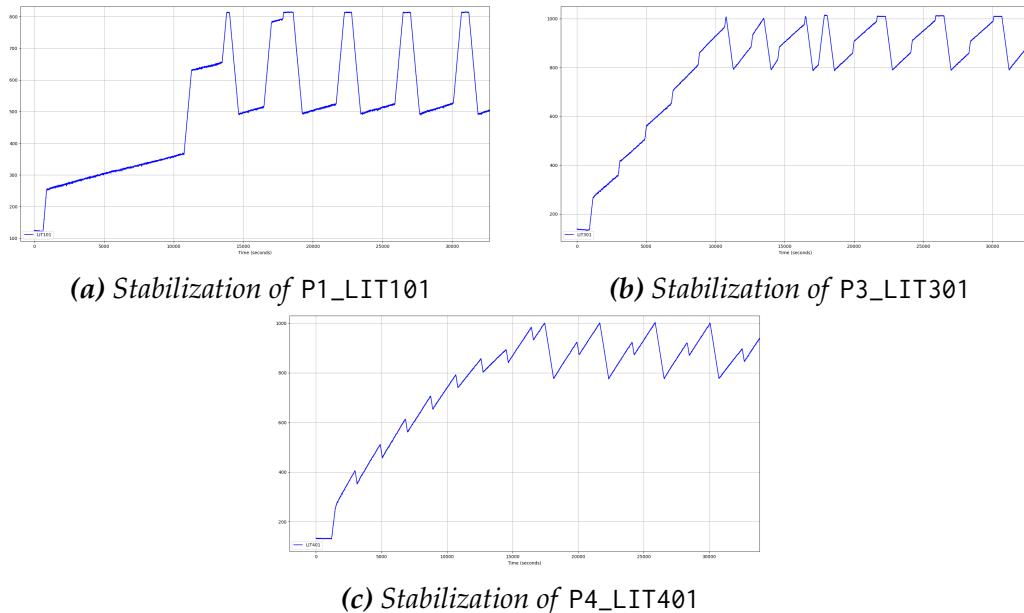


Figure 5.4: SWaT stabilization

2822 In total, the dataset consists of thousands of lines collecting 51 attributes
 2823 denoting the state of the system.

2824 **Network Traffic** The network traffic was collected using an appliance
 2825 from a well-known network hardware manufacturer and was made avail-
 2826 able only in CSV format and not PCAP format. Table 5.2 shows some of
 2827 the main data captured:

Category	Description
Date	Date of Log
Time	Time of Log
Origin	IP of server
Source IP	IP address of source
Destination IP	IP address of destination
Protocol	Network protocol
Application Name	Name of application
Modbus Function Code	Function Code
Modbus Function Description	Description of function
Modbus Transaction ID	Transaction ID
SCADA Tag	Sensor or actuator ID
Modbus Value	Value transmitted
Service / Destination Port	Port number of destination IP
Source Port	Port number of source IP

Table 5.2: SWaT network traffic data

2828 It is important to note that the presence of the string *Modbus* within the
 2829 name of the captured data fields might be misleading, as it may give the
 2830 impression that the communication is taking place via the Modbus proto-
 2831 col rather than the CIP over EtherNet/IP protocol. However, in reality, the
 2832 latter is the actual protocol being used. The indication of Modbus within
 2833 the network traffic dataset could be a result of the appliance used for net-
 2834 work data sniffing, which might have encapsulated CIP over EtherNet/IP
 2835 protocol packets within Modbus frames. This technique is described by
 2836 Snide in [74].

2837 To confirm that the communication protocol used is CIP over Ether-
 2838 Net/IP, we can examine the "Service / Destination Port" field. The destina-
 2839 tion port displayed is TCP port 44818, which corresponds to the TCP port
 2840 commonly used for transporting **explicit messages** in the EtherNet/IP
 2841 protocol (as explained in Section 2.2.6.2). Further evidence supporting the
 2842 use of CIP protocol is observed in the *Modbus Function Code* field, con-
 2843 sistently displaying a value of 76. In the Modbus protocol, there is no

2844 corresponding function associated with that code (as described in Section
2845 2.2.6.1). However, when converting the value 76 from decimal to hex-
2846 adecimal, it translates to **4C**, which corresponds to the **Read Tag Request**
2847 **service in the CIP protocol**. This correlation is further supported by the
2848 presence of the *Application Name* and *Modbus Function Description* fields,
2849 explicitly indicating the protocol name and service within them.

2850 Datasets for years beyond 2015 were not considered due to significant
2851 limitations that impede their usability. For instance, the 2017 dataset ex-
2852 hibits a high level of granularity in the physical system data, combined
2853 with short temporal periods for data gathering. As a result, valuable in-
2854 formation is lost, and the network traffic data cannot be effectively uti-
2855 lized. On the other hand, post-2017 datasets suffer from the issue of being
2856 "spurious." This means that no distinction is made between the data col-
2857 lected during normal operations of the SWaT system and the data associ-
2858 ated with system attacks, as was done in the 2015 dataset. Furthermore,
2859 the 2018 dataset lacks network traffic data altogether.

2860 Regarding network traffic related to the year 2017, the only ones that
2861 seemingly remain unaffected by attacks, the data is divided into large
2862 PCAP files weighing more than 6 GB each. Despite their size, these files
2863 only contain a few minutes of data, which is insufficient to cover even
2864 one complete system cycle. This renders the use of these files impractical,
2865 considering the available resources.

2866 Despite their large quantity, the CSV files associated with network traf-
2867 fic for the year 2015 are comparatively smaller in size, just slightly over 100
2868 MB each. These files cover a more extensive time period, which facilitates
2869 easier management and analysis. Prioritizing the physical process dataset
2870 as crucial, I have selected the year 2015 as a case study, even though the
2871 network data may be incomplete.

Our Framework at Work: Reverse Engineering of the iTrust SWaT System

2872 IN THIS chapter, our main objective is to apply the framework and method-
2873 ology introduced in Chapter 4 to the case study of the iTrust SWaT sys-
2874 tem, as illustrated in Chapter 5. The purpose of this analysis is to assess
2875 the effectiveness and potential of the proposed framework within the con-
2876 text of a system that closely replicates a real-world water treatment plant,
2877 albeit on a smaller scale.

2878 Due to the complexity of the system and the limited space available
2879 in this thesis, we will not conduct a comprehensive analysis and reverse
2880 engineering of the entire system. Instead, we will focus on specific parts
2881 for analysis. We leave it to the reader or those interested in utilizing the
2882 proposed methodology and framework to complete the analysis, should
2883 they choose to do so.

2884 By focusing on selective components and leaving room for further ex-
2885 ploration, we strike a balance between providing valuable insights and
2886 acknowledging the potential for additional research. This approach em-
2887 powers the reader and interested individuals to explore the iTrust SWaT
2888 system further and leverage the proposed methodology and framework
2889 for a more comprehensive analysis.

2890 6.1 Preliminary Operations

2891 Prior to beginning the actual analysis, several preliminary manual op-
2892 erations need to be conducted on the physical process dataset utilized as
2893 a case study, specifically the SWaT system dataset for the year 2015 as out-
2894 lined in Section 5.2.1. To simulate the data-capture process performed by
2895 Ceccato et al. using their scanning tool, the original dataset in XLSX format
2896 (proprietary to Microsoft Excel) was divided into multiple datasets in CSV
2897 format. Each of these datasets corresponds to the individual stages of the
2898 SWaT system and contains the respective registers. These resulting files
2899 were then saved in the directory specified by the `raw_dataset_directory`
2900 directive in the framework configuration file, `config.ini`, ready to be used
2901 in the pre-processing phase. Furthermore, the headers were manually re-
2902 named by adding a prefix from P1_ to P6_ to each register's name. This
2903 prefix indicates the stages, ensuring that each register is easily identifiable
2904 and linked to its corresponding stage.

2905 6.2 Planning the Analysis Strategy

2906 The complexity of the system being analyzed necessitates the adoption
2907 of a deliberate strategy for the analysis. It is not feasible to rely on trial and
2908 error or attempt every possible combination between stages. The former
2909 approach may overlook crucial relationships between PLCs or between
2910 registers, while the latter may result in excessive and unproductive efforts
2911 if the specific portion of the system being analyzed lacks significant infor-
2912 mation or relationships. A sound analysis strategy helps us focus on the
2913 important parts of the system, improving the quality of the analysis and
2914 leading to better process comprehension. By prioritizing our attention, we
2915 can gain a deeper understanding of the crucial components, resulting in
2916 more informed decision-making and a comprehensive understanding of
2917 the overall processes.

2918 To define this strategy, a potential starting point could involve analyz-

ing network traffic to determine the communication patterns and participants within the system. This can be accomplished by utilizing the techniques discussed in Section 4.2.2 on Network Analysis. By applying the Python script described in that section to the data extracted from the network traffic dataset debated in Section 5.2.1, we can generate a (simplified) network graph, as illustrated in Figure 6.1.

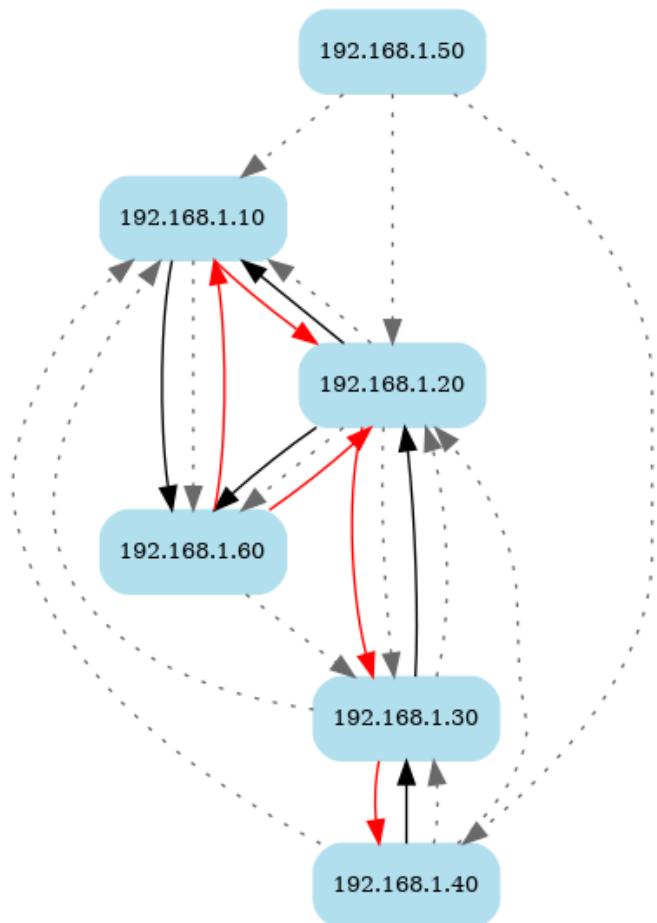


Figure 6.1: Simplified graph of the iTrust SWaT system network

The graph clearly illustrates the structure of communications between the PLCs. Referring back to Table 5.1, which displays the IP address - PLC associations, we can observe that PLCs 1 through 4 communicate directly and sequentially with each other in a Request/Response communication pattern (represented by red and black arrows, respectively). Additionally,

2930 PLC6 communicates with both PLC1 and PLC2. On the other hand, the
2931 gray dotted arrows indicate communications for which we have knowl-
2932 edge of a response, but the corresponding request is unknown. For the
2933 purposes of our analysis strategy, we will not consider these communica-
2934 tions within this context.

2935 Based on our observations, the analysis strategy we will adopt involves
2936 considering sequential pairs of PLCs to effectively capture the relation-
2937 ships and implications between registers.

2938 **6.3 Reverse Engineering of the iTrust SWaT Sys-
2939 tem**

2940 Before we delve into the analysis, it is important to provide some pre-
2941 liminary remarks.

2942 **Analysis structure** Firstly, the analysis will be structured as a schematic
2943 analysis due to space constraints, which prevent us from presenting the
2944 extensive inferences and reasoning regarding the system in full detail.
2945 Therefore, following the analysis strategy outlined in Section 6.2, we will
2946 concentrate exclusively on the pairs of PLCs comprising PLC1-2, PLC2-
2947 3, and PLC3-4. However, the general procedure of the methodology and
2948 how to reason about the data obtained from the framework have already
2949 been demonstrated through examples on the PLC1 of the SWaT system in
2950 Chapter 4. We encourage readers to refer to those examples for a more
2951 comprehensive understanding. In this analysis, our focus will be on illus-
2952 trating the conjectures and properties that arise from the analysis, utilizing
2953 tables and the outputs generated during the analysis.

2954 **Defining subsystem time duration** The second premise addresses the
2955 process of defining the subsystems to be analyzed, which were obtained
2956 during the pre-processing phase, during the merge phase of the individual

2957 datasets. Apart from the projection determined by the considered PLCs,
2958 a time-based selection of the analysis period has also been performed (see
2959 Section 4.2.3.1). This selection spans a duration of 20,000 seconds, which
2960 is equivalent to approximately five and a half hours or roughly five sys-
2961 tem cycles. The analysis begins at 100,000 seconds, which corresponds to
2962 approximately 27 hours from the start of the available data. This deliber-
2963 ate selection aims to exclude the initial transient period during which the
2964 SWaT system is initialized. We believe that this time range is more than
2965 sufficient for accurately defining the characteristics of the SWaT system
2966 components.

2967 **Conventions** The third premise introduces a convention that governs
2968 the naming of PLC registers and will be consistently followed throughout
2969 our analysis. According to this convention, registers with similar names,
2970 such as P1_LIT101 and P3_LIT301, P2_MV201 and P3_MV202, are considered
2971 to belong to the **same category or type of register**. This convention allows
2972 us to establish a relationship or correspondence between registers based
2973 on their naming pattern. By grouping registers with similar names, we can
2974 infer that they serve similar functions or represent similar components in
2975 the system, such as level sensors, tanks, pumps, and so on.

2976 **About the Business Process Analysis** In the end, the Business Process
2977 Analysis will focus solely on the physical process part. This is because the
2978 datasets of network traffic captures provided by iTrust for the year 2015
2979 (as discussed in Section 5.2.1) are **incomplete**. While communications re-
2980 lated to measurements are present, those associated with actuators are en-
2981 tirely missing, as well as additional communications related to other sys-
2982 tem characteristics that we observed in the datasets of subsequent years.
2983 As a result, we were unable to incorporate the network event recognition
2984 component into our Business Process Analysis. To implement this com-
2985 ponent, we would require complete and overlapping network data, along
2986 with a clean physical process dataset not affected by system attacks. Un-
2987 fortunately, none of the available iTrust datasets fulfill these criteria.

2988 6.3.1 Reverse Engineering of PLC1 and PLC2

2989 The initial focus of analysis will be on the pair comprising PLC1 and
 2990 PLC2. Let's delve into the main features of this subsystem by examining
 2991 the outcomes obtained from applying the framework to it.

2992 6.3.1.1 Pre-processing - Preliminary Analysis

2993 **Measurements and Actuators Recognition** Listing 6.1 shows the out-
 2994 comes obtained from automatic recognition of likely measurements and
 2995 actuators:

```

2996   1  Actuators:
2997   2    P1_MV101 [0.0, 1.0, 2.0]
2998   3    P1_P101 [1.0, 2.0]
2999   4    P2_MV201 [0.0, 1.0, 2.0]
3000   5    P2_P203 [1.0, 2.0]
3001   6    P2_P205 [1.0, 2.0]
3002   7
3003   8  Sensors:
3004   9    P1_FIT101 {'max_lvl': 2.7, 'min_lvl': 0.0}
3005  10    P1_LIT101 {'max_lvl': 815.1, 'min_lvl': 489.6}
3006  11    P2_AIT201 {'max_lvl': 256.5, 'min_lvl': 252.9}
3007  12    P2_AIT202 {'max_lvl': 8.4, 'min_lvl': 8.3}
3008  13    P2_AIT203 {'max_lvl': 342.8, 'min_lvl': 320.0}
3009  14    P2_FIT201 {'max_lvl': 2.5, 'min_lvl': 0.0}
3010  15
3011  16  Hardcoded setpoints or spare actuators:
3012  17    P1_P102 [1.0]
3013  18    P2_P201 [1.0]
3014  19    P2_P202 [1.0]
3015  20    P2_P204 [1.0]
3016  21    P2_P206 [1.0]
```

Listing 6.1: Preliminary analysis outcomes for sensors and actuators of PLC1–2

3017 Based on the results presented in Listing 6.1, the framework has iden-
 3018 tified P1_MV101, P1_P101, P2_MV201, P2_P203, and P2_P205 as **probable ac-**
 3019 **tuators**. The actuators denoted by the *Pxxx* notation are binary actuators

3020 (not boolean!), meaning they have two states represented by the values
3021 1 and 2. Conversely, the actuators identified by the *MVxxx* notation are
3022 ternary actuators with three distinct states: 0, 1, and 2.

3023 To simplify the analysis, we have arbitrarily categorized the registers
3024 identified by the notation *MVxxx* as **valves** and the registers identified by
3025 the notation *Pxxx* as **pumps**. It is important to note that this distinction
3026 is solely for convenience and does *not* necessarily reflect the actual role or
3027 function of these actuators within the system.

3028 P1_FIT101, P1_LIT101, P2_AIT201, P2_AIT202, P2_AIT203, and P2_FIT201
3029 have been identified as **likely measurements**. Upon analyzing the range
3030 of values for register P1_LIT101, we observe a significant difference be-
3031 tween the maximum and minimum values. This observation leads us
3032 to speculate that P1_LIT101 could be identified as a **level sensor for the**
3033 **tank controlled by PLC1**. However, when examining registers P1_FIT101,
3034 P2_FIT201, P2_AIT201, and P2_AIT202, the small difference between their
3035 maximum and minimum values makes it unlikely that they represent ad-
3036 dditional tanks.

3037 Regarding register P2_AIT203, although the range of values is not as
3038 wide as in the case of P1_LIT101, it is still worth examining more closely. It
3039 is possible that P2_AIT203 indicates the presence of a small tank. However,
3040 considering our speculation that the other P2_AIT20x registers are not tank
3041 level sensors, it is uncertain whether P2_AIT203 falls into that category as
3042 well. Further analysis is required to confirm its role within the system.

3043 Some registers have been identified as **hardcoded setpoints or spare**
3044 **actuators** based on their constant values. These registers exhibit similar-
3045 ties to the previously recognized pump registers. It is plausible to spec-
3046 ulate that these registers could correspond to **spare actuators**. Moreover,
3047 the constant value of 1 associated with these registers suggests that it may
3048 represent the **OFF state** of the pumps.

3049 **Actuator State Durations** To gain a deeper understanding of the differ-
3050 ent states (0, 1, and 2) associated with valves P1_MV101 and P2_MV201, we

3051 can analyze the duration of each state. Listing 6.2 provides information
 3052 regarding the duration (in seconds) of states for these specific actuators:

```

3053 1 Actuator state durations:
3054 2 P1_MV101 == 0.0
3055 3 9 9 10 9 9 10 9 9 10 9
3056 4
3057 5 P1_MV101 == 1.0
3058 6 1174 1168 1182 1160 1172
3059 7
3060 8 P1_MV101 == 2.0
3061 9 669 3019 3012 3000 2981
3062 10
3063 11 P2_MV201 == 0.0
3064 12 8 8 8 9 9 8 9 9 9 9
3065 13
3066 14 P2_MV201 == 1.0
3067 15 1057 1057 1045 1038 1039
3068 16
3069 17 P2_MV201 == 2.0
3070 18 120 3135 3144 3127 3109

```

Listing 6.2: Time duration of the states of actuators P1_MV101 and P1_MV201 of PLC1-2

3071 It is evident that the duration of **state 0** is relatively short, averaging
 3072 around 8-10 seconds, while the other states have much longer durations.
 3073 This observation suggests that state 0 of a valve is a **transient state**, in-
 3074 dicating a transitional phase within the valve cycle. However, without
 3075 further information, it is currently not possible to determine the specific
 3076 position of state 0 within the overall valve cycle.

3077 **Actuator State Changes** Now that we have identified P1_LIT101 as the
 3078 supposed level sensor of the tank, we can examine the trend of the tank
 3079 level as the actuators change state. Listing 6.3 provides information on
 3080 the levels of the tank in correlation with the state changes of the P1_P101
 3081 pump:

```

3082 1 Actuator state changes:
3083 2 ...

```

3084	3	P1_LIT101	P1_P101	prev_P1_P101
3085	4	536.0356	1	2
3086	5	533.3272	1	2
3087	6	542.1591	1	2
3088	7	534.8581	1	2
3089	8	540.5890	1	2
3090	9
3091	10	P1_LIT101	P1_P101	prev_P1_P101
3092	11	813.0031	2	1
3093	12	813.0031	2	1
3094	13	811.8256	2	1
3095	14	812.7283	2	1
3096	15	813.3171	2	1
3097	16

Listing 6.3: P1_P101 state changes in relation to P1_LIT101

3098 Based on the speculation that state 1 represents the OFF state of the
 3099 pump and state 2 represents the ON state, we can analyze the data in List-
 3100 ing 6.3. When pump P1_P101 switches from the ON state to the OFF state,
 3101 the average level of P1_LIT101 is 535. On the other hand, when P1_P101
 3102 goes from the OFF state to the ON state, the average level of P1_LIT101
 3103 is 813. These values correspond to the **minimum and maximum relative**
 3104 **setpoints** of P1_P101, respectively.

3105 Based on this information, we can infer that pump P1_P101 is respon-
 3106 sible for **emptying the tank**. Moreover, it can be extended to assume that
 3107 a pump, in general, is responsible for **water outflow**.

3108 Applying the same analysis to the data for valve P1_MV101, which is not
 3109 reported for conciseness, we can speculate that P1_MV101 is responsible for
 3110 **filling the tank**. In this case, states 1 and 2 would represent the **OFF and**
 3111 **ON states of the valve**, respectively. The relative setpoints of P1_MV101
 3112 are approximately 500 (minimum) and 800 (maximum). By extending this
 3113 analysis, we can speculate that a valve, such as P1_MV101, is responsible
 3114 for controlling the **water inflow**.

3115 Regarding the elements controlled by PLC2 and the sensor P1_FIT101,

3116 the analysis does not reveal the presence of another tank. Therefore, we
 3117 cannot determine the exact role of sensors P2_AIT20x, P1_FIT101 and P2_FIT201
 3118 at this point. However, there is a similarity observed between the relative
 3119 setpoints of P1_P101 and those of P2_MV201, P2_P203, and P2_P205. These
 3120 registers exhibit very similar values during state changes, suggesting a
 3121 potential relationship or similar control behavior between them.

3122 6.3.1.2 Graphs and Statistical Analysis

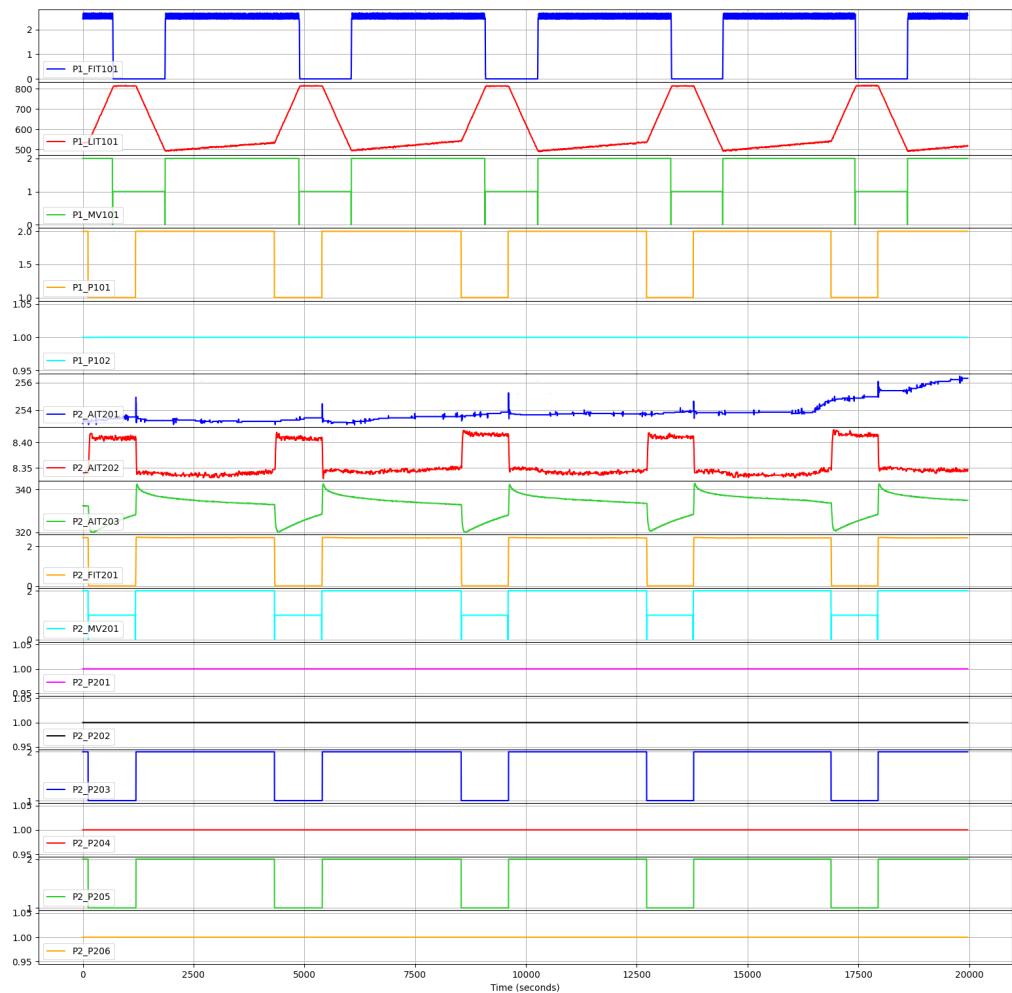


Figure 6.2: Chart of PLC1-2 registers

3123 Figure 6.2 illustrates the graphical representation of the registers in

3124 PLC1 and PLC2 and their respective trends.

3125 The image provides additional support to the conjectures made during
 3126 the preliminary analysis regarding the spare actuators. Furthermore, it is
 3127 evident from the graphs that these spare actuators **do not appear to influ-**
 3128 **ence the trend** of any of the measurements. Therefore, based on this obser-
 3129 vation, we can confidently exclude these registers from further graphical
 3130 analysis.

3131 Figure 6.3 shows a clearer representation of the subsystem after remov-
 3132 ing the spare actuators.

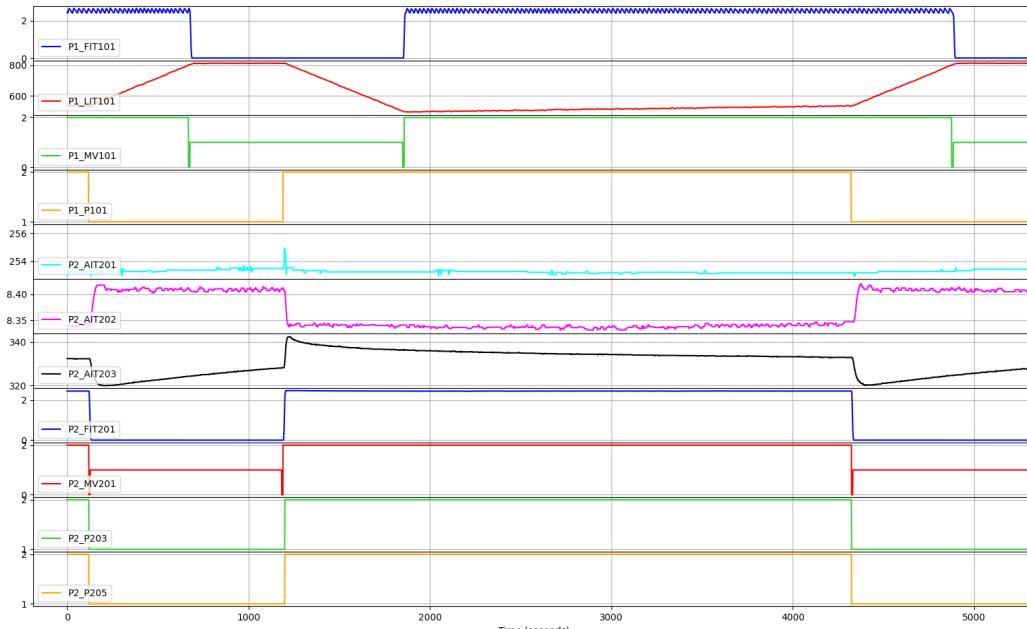


Figure 6.3: Chart of PLC1-2 registers without spare actuators (particular)

3133 Figure provides furthermore additional insights that allow us to spec-
 3134 ulate on aspects that remained unexplained during the preliminary analy-
 3135 sis. The most prominent aspect that stands out is the relationship between
 3136 the level behavior of P1_LIT101 and the states of P1_P101 and P1_MV101. It
 3137 appears evident that these two actuators **do not exhibit complementary**
 3138 **behavior**, meaning that their states do not alternate in an ON-OFF pattern.
 3139 Instead, they can remain in the same state, either ON or OFF, for extended

3140 periods of time. When they are both in the ON state, **the growth of the**
3141 **tank level is slow**. However, as soon as P1_P101 switches to the OFF state,
3142 the tank level starts to **increase at a faster rate**. Therefore, when both of
3143 these actuators are in the ON state the influx of water into the tank ex-
3144 ceeds the outflow of water from it. Conversely, when P1_MV101 switches
3145 to the OFF state, the tank level **decreases**. During periods when both ac-
3146 tuators are in the OFF state, there is a *plateau* where the tank level remains
3147 relatively **stable**.

3148 Furthermore, we observe a relationship between P1_FIT101 and the
3149 trend of P1_MV101. When the valve is in the off state, P1_FIT101 registers
3150 a value of 0, whereas it registers a value greater than 2 when the valve is
3151 open. This suggests that P1_FIT101 could be a **sensor associated with the**
3152 **flow** of water entering the tank, which is represented by P1_LIT101. By
3153 drawing an analogy with its name, it is plausible to consider P2_FIT201 as
3154 another flow sensor.

3155 Another intriguing aspect that arises from examining the graphs in Fig-
3156 ure 6.2 and Figure 6.3 is the **non-cyclic pattern** observed in the P2_AIT201
3157 measurement. Instead of exhibiting a cyclical trend, it follows a linear
3158 pattern. Furthermore, considering the narrow range of values associated
3159 with this measurement, it is reasonable to speculate that P2_AIT201 may be
3160 associated with a **sensor that measures a specific property of the water**.

3161 The limited range of values observed in P2_AIT202 also raises the pos-
3162 sibility that it functions as a sensor for a particular water characteristic,
3163 despite exhibiting a cyclic pattern. As for P2_AIT203, its role remains un-
3164 defined, although it also displays a cyclic trend. This trend, along with that
3165 of P2_AIT202, does not appear to be related to the behavior of the valves
3166 (which rules out the possibility of it being related to a tank), but rather
3167 to that of the pumps. Consequently, it is imperative to conduct further
3168 investigations into these aspects.

3169 By examining the trend of valve P2_MV201, an additional speculation
3170 can be made. It appears to be independent of the trends observed in any

of the measurements within this subsystem. Based on previous conjectures regarding the role of valves and considering the duration of its ON and OFF states, it is possible that P2_MV201 is responsible for **filling a tank that is not part of this particular subsystem**. Once again, a thorough investigation is necessary to confirm this hypothesis.

6.3.1.3 Invariant Inference and Analysis

Through the process of *invariant analysis*, we aim to discover new information about the system and determine whether the conjectures made in the previous steps are supported by the data obtained from the Daikon analysis.

General Invariants We will begin this phase by analyzing the general invariants (see Section 4.2.5.1.). Listing 6.4 presents a selection of these invariants:

```
3184 1 P2_P206 == P2_P204 == P2_P202 == P2_P201 == P1_P102 == 1.0
3185 2 P2_P205 == P2_P203
3186 3 max_P1_LIT101 == 816.0
3187 4 min_P1_LIT101 == 489.0
3188 5 max_P2_AIT201 == 257.0
3189 6 min_P2_AIT201 == 252.0
3190 7 max_P2_AIT202 == 9.0
3191 8 min_P2_AIT202 == 8.0
3192 9 max_P2_AIT203 == 343.0
3193 10 min_P2_AIT203 == 320.0
```

Listing 6.4: General Invariants for PLC1-2

The invariant mentioned on line 2 is particularly significant: it states that P2_P203 and P2_P205 always have the same values. While this information was somewhat apparent in the previous steps, it becomes more apparent and evident in this analysis. This observation leads us to speculate that the two pumps, P2_P203 and P2_P205, **are related to each other** in some way. The other invariants provided in this section further reinforce the hypotheses about the spare actuators.

3201 **Analysis on Single Actuator States** We proceed with the examination
 3202 of the invariants derived from the *first of the two semi-automatic analysis*
 3203 discussed in Section 4.2.5.2. Specifically, we will focus on the analysis con-
 3204 cerning the states of individual actuators in relation to a specific measure-
 3205 ment, which in our case pertains to the tank represented by P1_LIT101. For
 3206 illustrative purposes, let's consider states 1 and 2 (OFF e ON respectively)
 3207 of valve P1_MV101 as an example (we will disregard state 0 as it is consid-
 3208 ered transient). The conditional invariants pertinent to this scenario can
 3209 be found in Listing 6.5.

```

3210 1 =====
3211 2 P1_MV101 == 1.0 && P1_LIT101 < max_P1_LIT101 - 16 &&
3212   ↪ P1_LIT101 > min_P1_LIT101 + 15
3213 3 =====
3214 4 ...
3215 5 P2_P205 == P2_P203 == P2_MV201 == P1_P101 == 2.0
3216 6 P1_FIT101 == 0.0
3217 7 slope_P1_LIT101 == -1.0
3218 8 P2_FIT201 > P1_FIT101
3219 9 P2_FIT201 > P1_MV101
3220 10 P2_FIT201 > P1_P101
3221 11 ...
3222 12
3223 13 =====
3224 14 P1_MV101 == 2.0 && P1_LIT101 < max_P1_LIT101 - 16 &&
3225   ↪ P1_LIT101 > min_P1_LIT101 + 15
3226 15 =====
3227 16 slope_P1_LIT101 == P1_P102
3228 17 P1_FIT101 > P1_MV101
3229 18 ...
3230 19 P1_MV101 >= P1_P101
3231 20 P1_MV101 >= P2_MV201
3232 21 P1_MV101 >= P2_P203
3233 22 P2_P203 >= P1_P101
3234 23 ...

```

Listing 6.5: Conditional Invariants for states 1 and 2 of P1_MV101

3235 To prevent transient periods caused by water flow stabilization when

3236 the actuators change state, a condition is imposed on the level of P1_LIT101.
3237 However, this condition may result in an incomplete understanding of the
3238 system's behavior. To address this, a manual refinement of the analysis is
3239 required, utilizing the runDaikon.py script as outlined in Section 4.2.5.2.

3240 Based on the analysis, the following observations can be made when
3241 the valve is in the OFF state:

- 3242 • The slope of P1_LIT101, denoted as slope_P1_LIT101, is negative
3243 (line 7). This indicates a **downward trend** in the tank level, as we
3244 have seen in Section 6.3.1.1.
- 3245 • P1_P101 is in state 2, or ON (line 5).
- 3246 • P1_FIT101 is zero (line 6).
- 3247 • P2_FIT201 has a value greater than 2 (line 10).

3248 On the other hand, when the valve is in the ON state:

- 3249 • slope_P1_LIT101 is positive (line 16). This indicates an **upward trend**
3250 in the tank level, as we have seen in Section 6.3.1.1.
- 3251 • P1_FIT101 assumes a value greater than 2. The combination of this
3252 finding, along with the previous one regarding the same register,
3253 strengthens the hypothesis that this is indeed a flow sensor.
- 3254 • P1_P101 can be in either the ON or OFF state, as we have seen in
3255 Section 6.3.2.2.

3256 When conducting a manual analysis using the runDaikon.py script on
3257 tank levels that fall outside the range defined by the previous condition, it
3258 does not yield useful slope information. This situation can occur because,
3259 despite the noise attenuation applied to the tank level sensor data, if there
3260 is even a single cycle in the system where the calculated slope deviates
3261 from the expected outcome, it can adversely affect the entire Daikon anal-
3262 ysis. Figure 6.4 shows this behavior.

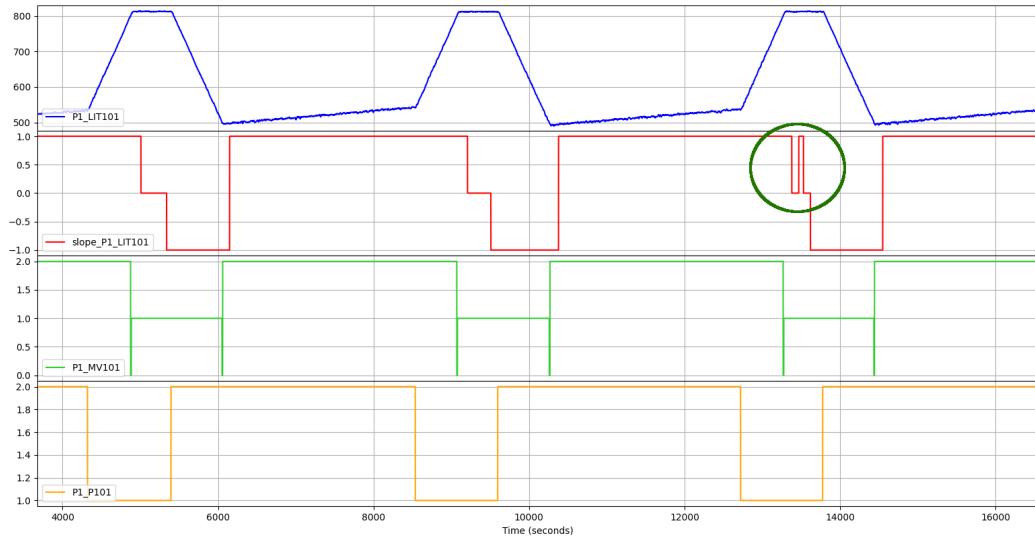


Figure 6.4: Slope calculation anomaly (in the circle)

3263 **Analysis of the Current System Configuration** We conclude our analy-
 3264 sis of invariants by considering the *second semi-automatic analysis* outlined
 3265 in Section 4.2.5.2, which focuses on the effective states of the system. Due
 3266 to the comprehensive nature of this analysis, we will not provide a de-
 3267 tailed report of the outputs to maintain brevity. However, this analy-
 3268 sis confirms the findings observed in the analysis of individual actuator
 3269 states. Additionally, one crucial piece of information becomes apparent
 3270 for future steps: the changing state of the actuators controlled by PLC2 **do**
 3271 **not impact the behavior of the tank** controlled by PLC1. Indeed, it is suf-
 3272 ficient to examine the invariant pertaining to the slope of the tank to verify
 3273 this observation.

3274 6.3.1.4 Business Process Mining and Analysis

3275 As explained in Section 4.2.6.1, the *process mining phase* applied to the
 3276 physical system provides us with an immediate understanding of the sys-
 3277 tem cycle and the chronological sequence of states. It enables us to de-
 3278 termine the duration of time the system remains in a particular state and
 3279 at what relative setpoint the state transition occurs concerning the refer-

ence measure. Furthermore, we can analyze the trend of this measure within each state. Additionally, we can examine the relative setpoints of other measurements to identify any connections between changes in system state and these values.

Given that we have already identified the likely measurement representing the tank and the corresponding actuators that control its behavior, an activity diagram can be generated as depicted in Figure 6.5.

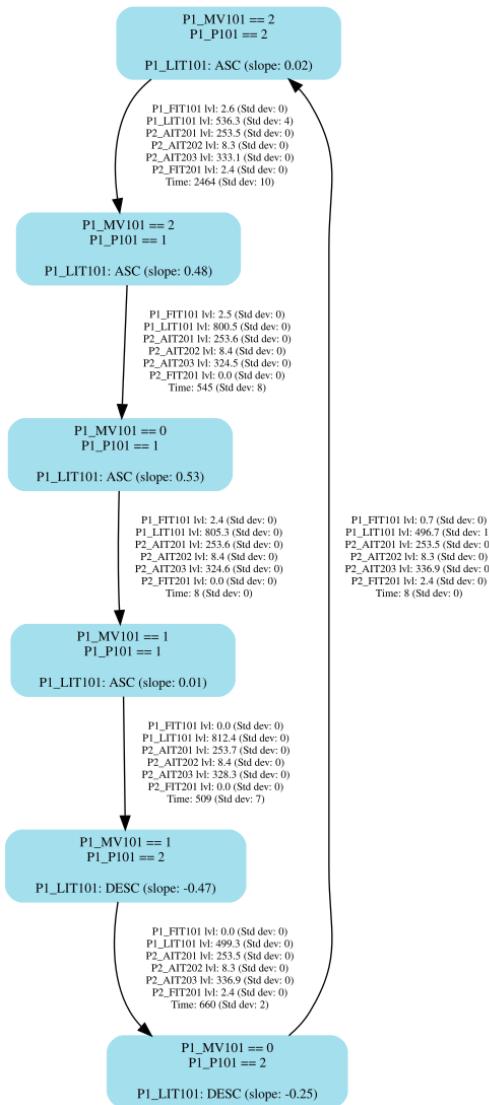


Figure 6.5: Activity diagram for PLC1-2

3287 The activity diagram allows for easy interpretation of the tank level
3288 trend and slope within different states. The states where the valve has a
3289 value of 0 can be disregarded due to their short duration. Similar to the
3290 graphical analysis, there is an observed change in slope during the increas-
3291 ing trend of the tank level between the states [P1_MV101 == 2, P1_P101
3292 == 2] and [P1_MV101 == 2, P1_P101 == 1], where the slope changes from
3293 0.02 to 0.48.

3294 Additionally, the timing analysis on the edges reveals that the tank
3295 takes longer to fill than to empty, and the system remains in the [P1_MV101
3296 == 1, P1_P101 == 1] state for approximately 8 minutes (509 seconds).

3297 Regarding the state [P1_MV101 == 1, P1_P101 == 1], there appears
3298 to be a discrepancy between the trend of the tank level as reported in the
3299 activity diagram and the conjectures made in the previous phases of the
3300 analysis. The activity diagram correctly depicts an increasing trend in the
3301 tank level between the end of the state [P1_MV101 == 2, P1_P101 == 1]
3302 and the end of the state [P1_MV101 == 1, P1_P101 == 1]. However, the
3303 discrepancy arises due to the fact that the interruption of flow during the
3304 valve state change from state ON to state OFF is not immediate, mainly
3305 due to the presence of the transient state 0. Additionally, water continues
3306 to flow within the piping towards the tank for a short duration even after
3307 the valve is closed. After this period, which usually lasts a few seconds,
3308 the tank level stabilizes, and there is no further inflow or outflow of water.

3309 By adjusting the tolerance parameter -t in the processMining.py script,
3310 it is possible to obtain accurate data regarding the behavior of the state
3311 corresponding to the *plateau* observed in the graphical analysis.

3312 The data presented on the arcs in the activity diagram represents the
3313 measurement values relative to the time of system state changes, specifi-
3314 cally the relative setpoints. These values are calculated based on the av-
3315 erage data collected in each cycle. By analyzing this data, we can observe
3316 the tank level values at which the system undergoes configuration changes
3317 and how the trend of the tank level changes accordingly. Specifically, we
3318 can see that the trend changes from ascending to stable at a tank level

3319 value of 800, from stable to descending at approximately 812, and from
3320 descending back to ascending at around 499. The change in the speed of
3321 tank filling occurs at approximately 535.

3322 Furthermore, the data provides additional support for the hypothesis
3323 that the measurements associated with registers P2_AIT20x are not influ-
3324 enced by the tank's trends.

3325 **6.3.1.5 Properties**

3326 From the conjectures derived from the four phases of the analysis we
3327 will derive the properties of the subsystem we are studying, which will be
3328 placed within a **summary table**. This table contains an integer identifying
3329 the property, the statement of the property itself, and from which of the
3330 four phases it was derived.

#	Statement	Derived from
P1	The registers P1_LIT101, P1_FIT101, P2_AIT201, P2_AIT202, P2_AIT203, and P2_FIT201 hold likely sensor measurements.	Preliminary Analysis Graphical Analysis
P2	The registers P1_MV101, P1_P101, P2_MV201, P2_P203, and P2_P205 holds likely actuator commands.	Preliminary Analysis Graphical Analysis
P3	The actuators that contain the substring "MVxxx" are considered to be three-state actuators. For simplicity, we refer to them as valves.	Preliminary Analysis
P4	The state 0 of these valves is associated to a transient state that occurs during the transition between state 1 (OFF) and state 2 (ON).	Preliminary Analysis Graphical Analysis
P5	The actuators that contain the substring "Pxxx" are considered to be binary actuators. For simplicity, we refer to them as pumps.	Preliminary Analysis
P6	The registers P1_P102, P2_P201, P2_P202, P2_P204, and P2_P206 are associated to spare actuators. They do <i>not</i> influence the trend of any measurements and they are considered in the OFF state.	Preliminary Analysis Graphical Analysis Invariant Analysis

P7	P1_LIT101 is associated to the level sensor of the tank controlled by PLC1.	Preliminary Analysis Graphical Analysis
P8	P1_MV101 and P1_P101 are the actuators responsible for the level behavior of the water contained in the tank.	Graphical Analysis Invariant Analysis Business Process
P9	P1_MV101 is responsible for filling the tank. P1_P101 is responsible for emptying the tank.	Graphical Analysis Invariant Analysis
P10	Valve are responsible for water inflow. Pumps responsible for water outflow.	Preliminary Analysis Graphical Analysis Invariant Analysis
P11	The rate of tank level growth is slow when both P1_MV101 and P1_P101 are in the ON state. The growth speed increases when P1_P101 switches to the OFF state.	Graphical Analysis Business Process
P12	When both P1_MV101 and P1_P101 are in the ON state, the inflow of water into the tank surpasses the outflow of water from the tank.	Graphical Analysis Business Process
P13	The tank level decreases when P1_MV101 is in the OFF state and P1_P101 is in the ON state.	Graphical Analysis Invariant Analysis
P14	The tank level remains (relatively) stable when both P1_MV101 and P1_P101 are in the OFF state.	Graphical Analysis Business Process
P15	The trend of the tank level changes from ascending to stable when the level reaches approximately 800. It shifts from stable to descending when the level averages around 812. It changes from descending back to ascending when the level reaches about 500. The speed of tank filling increases noticeably at around 535.	Business Process
P16	Absolute setpoints are 800, 812, 500 and 535.	Business Process
P17	P1_FIT101 serves as a flow or pressure sensor to measure the inflow or the pressure of water into the tank.	Graphical Analysis Invariant Analysis Business Process
P18	None of the actuators connected to PLC2 have an impact on the level of the tank controlled by PLC1.	Graphical Analysis Invariant Analysis Business Process
P19	None of the measurements connected to PLC2 represent a tank.	Preliminary Analysis Graphical Analysis
P20	P2_AIT201 does not exhibit a cyclic trend.	Graphical Analysis

P21	Both P2_AIT201 and P2_AIT202 serve as sensors for measuring certain properties of the water.	Preliminary Analysis Graphical Analysis
P22	The behavior and trend of P2_AIT202 and P2_AIT203 are directly associated with the operation of the pumps.	Graphical Analysis Invariant Analysis

Table 6.1: Properties of the PLC1-2 subsystem

6.3.2 Reverse Engineering of PLC2 and PLC3

Continuing our analysis of the iTrust SWaT system, our current focus will be on the registers of PLC3 and any potential relationships they may have with the registers of PLC2.

6.3.2.1 Pre-processing - Preliminary Analysis

Measurements and Actuators Recognition Listing 6.6 shows the outcomes obtained from automatic recognition of likely measurements and actuators. After previously identifying the measurements and actuators of PLC2 in Section 6.3.1.1, the listing exclusively showcases the registers associated with PLC3.

```

3341 1 Actuators:
3342 2 ...
3343 3 P3_MV301 [0.0, 1.0, 2.0]
3344 4 P3_MV302 [0.0, 1.0, 2.0]
3345 5 P3_MV303 [0.0, 1.0, 2.0]
3346 6 P3_MV304 [0.0, 1.0, 2.0]
3347 7 P3_P302 [1.0, 2.0]
3348 8
3349 9 Sensors:
3350 10 ...
3351 11 P3_DPIT301 {'max_lvl': 20.4, 'min_lvl': 0.0}
3352 12 P3_FIT301 {'max_lvl': 2.4, 'min_lvl': 0.0}
3353 13 P3_LIT301 {'max_lvl': 1014.5, 'min_lvl': 786.5}
3354 14
3355 15 Hardcoded setpoints or spare actuators:
3356 16 ...

```

3357 17 P3_P301 [1.0]

Listing 6.6: Preliminary analysis outcomes for sensors and actuators of PLC2-3

3358 From the provided listing, it is evident that the **likely measurements**
 3359 related to PLC3 are P3_DPIT301, P3_FIT301, and P3_LIT301. Drawing an
 3360 analogy with the derived properties from Table 6.1, P3_LIT301 can be as-
 3361 sociated with a **tank level sensor**, while P3_FIT301 may be linked to a
 3362 flow or pressure sensor. However, the specific role of P3_DPIT301 cannot
 3363 be speculated upon at this time.

3364 Regarding the **likely actuators**, they include P3_MV301, P3_MV302, P3_MV303,
 3365 P3_MV304, and P3_P302. By drawing parallels with the previous analysis
 3366 outcomes, it can be inferred that registers P3_MV30x represent valves, while
 3367 P3_P302 corresponds to a pump.

3368 Lastly, there is a **spare actuator** identified as P3_P301. Similar to the
 3369 previous analysis, this is an inactive pump, indicated by the constant value
 3370 of 1 in this register, signifying that it is in the OFF state.

3371 **Actuator State Durations** Let's proceed with the analysis of the actua-
 3372 tor states' duration, as displayed in Listing 6.7. In this analysis, our focus
 3373 will not be on examining the correspondence between values and the ac-
 3374 tual actuator states, as in Section 6.3.1.1. Instead, we will explore whether
 3375 these actuators exhibit any distinct patterns or behaviors based on their
 3376 duration.

```
3377 1 Actuator state durations:  

  3378 2 ...  

  3379 3 P3_MV301 == 1.0  

  3380 4 2527 4154 4154 4154 4094  

  3381 5  

  3382 6 P3_MV301 == 2.0  

  3383 7 36 35 36 35 34  

  3384 8  

  3385 9 P3_MV302 == 1.0  

  3386 10 662 138 654 138 656 139 658 137 656 137
```

```

3387 11
3388 12 P3_MV302 == 2.0
3389 13 62 1783 1596 1787 1591 1791 1576 1803 1540 1782
3390 14
3391 15 P3_MV303 == 1.0
3392 16 2526 4089 4088 4089 4028
3393 17
3394 18 P3_MV303 == 2.0
3395 19 97 96 97 96 96
3396 20
3397 21 P3_MV304 == 1.0
3398 22 689 1832 2206 1838 2203 1840 2191 1852 2152 1831
3399 23
3400 24 P3_MV304 == 2.0
3401 25 43 87 42 89 43 88 43 88 43 88
3402 26
3403 27 P3_P302 == 1.0
3404 28 637 115 632 115 632 114 634 114 632 115
3405 29
3406 30 P3_P302 == 2.0
3407 31 60 1821 1632 1825 1629 1829 1615 1841 1578 1820

```

Listing 6.7: Time duration of the states of actuators of PLC3

A notable behavior is observed in the P3_MV30x valves. P3_MV301, P3_MV303, and P3_MV304 have a relatively **short duration in the ON state**, ranging from around 30 seconds to a minute and a half. In contrast, P3_MV302 remains in the ON state for a longer duration but exhibits approximately twice as many cycles as the other actuators (10 cycles compared to 5 cycles). A similar characteristic is also observed in the OFF state of these actuators.

This behavior displayed by the actuators warrants further investigation in subsequent steps. However, based on the short duration of the ON state for P3_MV301, P3_MV303, and P3_MV304, **it appears unlikely that they have a significant impact on the tank level**. On the other hand, the influence of P3_MV302 cannot be ruled out and requires additional examination.

We can observed that P3_P302, the pump in PLC3, exhibits a behavior

similar to P3_MV302, with a number of cycles equal to 10. Furthermore, the durations of the ON and OFF states for both actuators appear to be overlapping. This suggests a **potential relationship between the two actuators**. Additionally, considering that P3_P302 is the only pump in PLC3, it is reasonable to speculate that it may have an influence on the tank level. Further investigation is necessary to validate this speculation and explore the precise nature of the relationship between P3_P302 and P3_MV302.

Actuator State Changes Based on the analysis of the probable measurements, we have identified the likely tank level sensor, represented by register P3_LIT301. In the previous analysis of the PLC1-2 subsystem in Section 6.3.1, the role of valve P2_MV201 remained unresolved. We speculated that this actuator might be responsible for the incoming water flow to an element outside the analyzed subsystem. To investigate this further, we can examine the relationship between P2_MV201 and the tank within this subsystem. By extracting information from the corresponding setpoints, we can gather insights to verify if our speculation holds true.

Listing 6.8 displays the setpoints associated with the state change of P2_MV201 in relation to the level of tank P3_LIT301.

```
3439 1 Actuator state changes:
3440 2 ...
3441 3 P2_MV201 prev_P2_MV201 P3_LIT301
3442 4 0 2 1000.2240
3443 5 0 1 799.1140
3444 6 0 2 1001.5060
3445 7 0 1 799.1942
3446 8 0 2 1001.5460
3447 9 0 1 799.1140
3448 10 ... ... ...
```

Listing 6.8: P2_MV201 state changes in relation to P3_LIT301

The setpoints provided in the listing support our conjecture. The maximum relative setpoint of 1000 and the minimum relative setpoint very close to 800 indicate a correlation that appears intentional. Based on this

3452 information, we can speculate that P2_MV201 is indeed the **valve responsible**
3453 **for filling the tank** associated with P3_LIT301.

3454 Regarding further information obtained from this step, there is limited
3455 insight available. While P3_P302 appears to be the pump responsible for
3456 emptying the tank associated with P3_LIT301, the analysis of the actua-
3457 tor lifetime indicates twice as many values compared to other actuators.
3458 The pump exhibits setpoint values of approximately 850 and 970 for the
3459 transition from ON to OFF, and 900 and 1000 for the transition from OFF
3460 to ON. On the other hand, P3_MV302 shows numerous state changes and
3461 shares setpoints with values close to 850, 970, 1000, and 900 for the same
3462 transitions. These values align perfectly with those of the P3_P302 pump,
3463 suggesting a potential relationship between the two actuators.

3464 Obtaining information about the remaining registers is challenging at
3465 this stage. Further analysis steps are required to gather additional insights
3466 and uncover more information about the system.

3467 6.3.2.2 Graphs and Statistical Analysis

3468 Graphical analysis can provide valuable insights and help validate the
3469 conjectures made during the preliminary analysis, as well as uncover con-
3470 nections between registers that were not identified in the previous step.
3471 To begin, we will test the hypothesis that valve P2_MV201 is responsible for
3472 filling the tank associated with P3_LIT101. Figure 6.6 displays these regis-
3473 ters, along with other registers whose behavior and relationships with the
3474 tank level we will explore in an attempt to gain a deeper understanding.

3475 Figure 6.6 shows the particular behavior of P3_LIT103, with two slope
3476 changes during the tank filling period. These slope changes correspond
3477 approximately to the relative setpoints found for P3_MV302 and P3_P302.
3478 However, we will analyze this aspect later. What we can see in relation
3479 to the initial conjecture about the role of P2_MV201 is that the period dur-
3480 ing which the valve remains in the ON state corresponds exactly to the
3481 increasing trend of P3_LIT301. The conjecture thus finds further support.

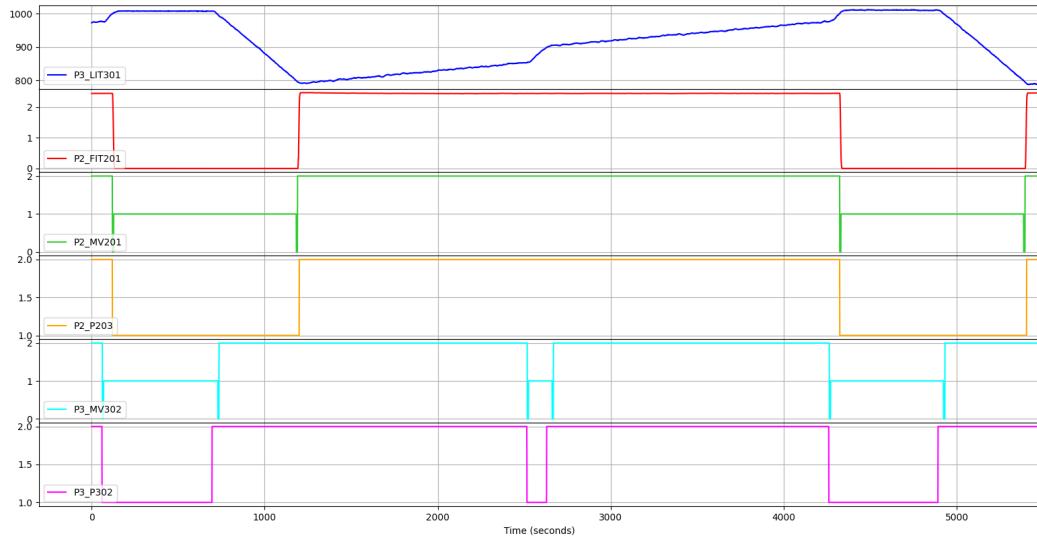


Figure 6.6: Verifying the conjecture about valve P2_MV201

3482 We also note how P2_FIT201 is related to the trend of P2_MV201 and the
 3483 increasing trend of P3_LIT301.

3484 Now let's examine the relationships between the tank represented by
 3485 P3_LIT301 and the other actuators of PLC3 through Figure 6.7.

3486 From the charts, it is evident that the trends of P3_DPIT301 and P3_FIT301
 3487 exhibit similarities. Furthermore, their overall pattern closely follows that
 3488 of the valve P3_MV302. Based on these observations, we can speculate that
 3489 there is a relationship between these registers or that they serve similar
 3490 functions, possibly as **pressure or flow sensors**.

3491 Regarding pump P3_P302, we observe that its OFF state coincides with
 3492 the increasing slope of P3_LIT301 during its upward trend and the entire
 3493 phase when the level remains relatively stable. Conversely, its ON state
 3494 corresponds to the gradual increase and decrease of the water level in the
 3495 tank. This observation provides further evidence to support the hypothe-
 3496 sis that P3_P302 is responsible for emptying the tank.

3497 Valve P3_MV302 exhibits a similar pattern to pump P3_P302. Building
 3498 upon our previous findings, we can speculate that, similar to P2_MV201 in
 3499 the previous analyzed subsystem, P3_MV302 is responsible for controlling

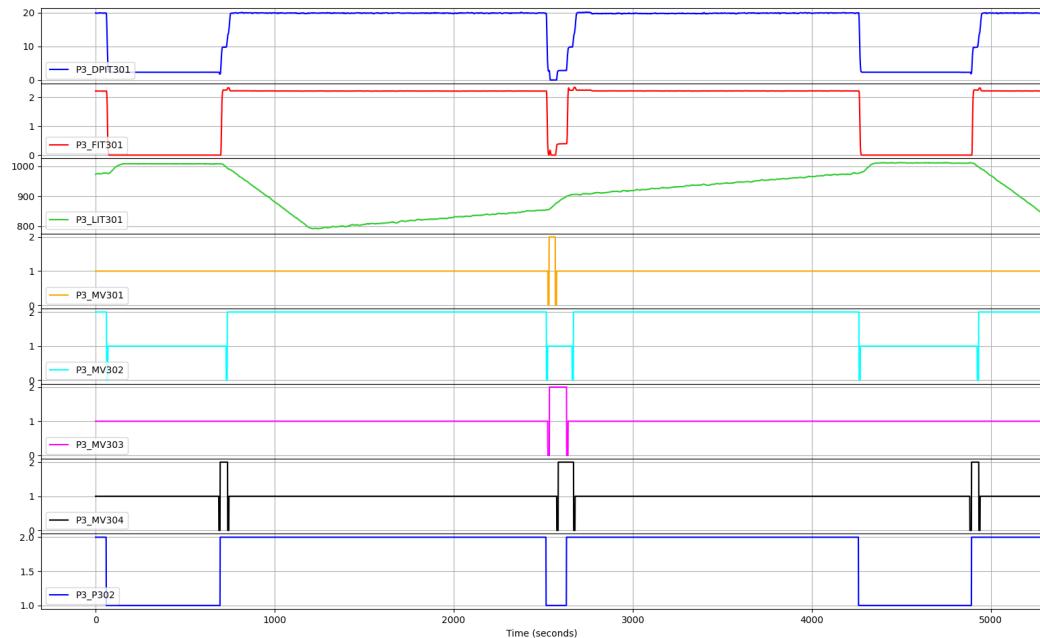


Figure 6.7: PLC3 registers

3500 the incoming flow to another element outside the analyzed subsystem.

3501 Let us now analyze the potential roles of valves P3_MV301, P3_MV303,
 3502 and P3_MV304 in relation to the indicated tank level. It seems unlikely
 3503 that P3_MV304 has any direct impact on the tank level. The valve is acti-
 3504 vated twice within one system cycle, but there are no noticeable changes
 3505 in P3_LIT301 during its first opening. This suggests that the second open-
 3506 ing is also insignificant in terms of tank level. However, it is worth noting
 3507 the slight peaks in P3_FIT301 that occur shortly after the valve's opening.

3508 Regarding the remaining valves, P3_MV301 and P3_MV303, their impact
 3509 on the tank level is still unclear, particularly during the increased slope of
 3510 P3_LIT101 around the second 2500. Figure 6.8 provides us with a compre-
 3511 hensive overview of the situation.

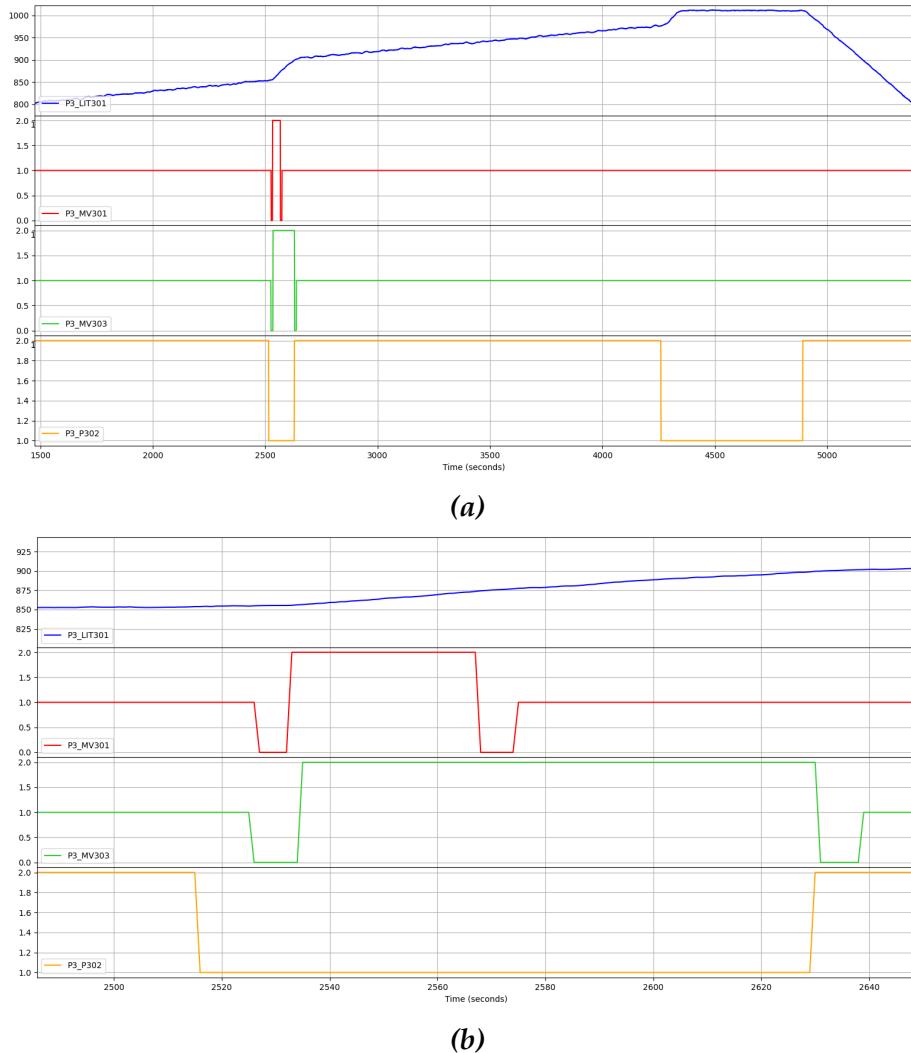


Figure 6.8: P3_MV301 and P3_MV303 analysis

3512 The analysis of the valves P3_MV301 and P3_MV303 indeed presents some
 3513 challenges. However, upon closer examination of Figure 6.8b, we observe
 3514 that when these valves are in the ON state, the slope of P3_LIT301 remains
 3515 relatively constant. This is unexpected, as one would anticipate a steeper
 3516 slope due to the inflow of liquid. Additionally, in Figure 6.8a, we can ob-
 3517 serve that when these valves are in the OFF state, the slope of P3_LIT301
 3518 from the second 4300 remains similar to the section we are currently ana-
 3519 lyzing. Based on these observations, we speculate that these valves either

3520 do not affect the water level of the tank or have a minimal, undetectable
 3521 impact.

3522 Indeed, Figure 6.7 provides additional insights into the relationship
 3523 between the valves P3_MV301, P3_MV303, and P3_MV304 and the sensors
 3524 P3_DPIT301 and P3_FIT301. The graph shows that the values of P3_DPIT301
 3525 and P3_FIT301 undergo noticeable changes during the activation of these
 3526 valves, suggesting a potential connection between them. This observation
 3527 supports the hypothesis that the valves and sensors are linked in some
 3528 way.

3529 **6.3.2.3 Invariant Inference and Analysis**

3530 **General Invariants** The analysis of general invariants does not yield any
 3531 significant insights. However, it confirms the maximum and minimum
 3532 values of the measurements seen in Section 6.3.2.1 and identifies the pres-
 3533 ence of the spare actuator P3_P301.

3534 **Analysis on Single Actuator States** The analysis of single actuator states
 3535 reveals additional information. The resulting invariants provide further
 3536 support for the conjecture regarding the roles of P2_MV201 and P3_P302 in
 3537 regulating the water level in the tank, with the former responsible for fill-
 3538 ing and the latter for emptying. Listing 6.9 presents the specific invariants
 3539 involved in this analysis.

```

3540 1 =====
3541 2 P2_MV201 == 1.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3542 3   ↪ P3_LIT301 > min_P3_LIT301 + 24
3543 4 =====
3544 5 P3_P301 == P3_MV303 == P3_MV301 == P2_P205 == P2_P203 ==
3545 6   ↪ P2_MV201 == 1.0
3546 7 P3_P302 == 2.0
3547 8 slope_P3_LIT301 == -1.0
3548 9 P3_FIT301 > P2_MV201
3549 10 P3_DPIT301 > P3_FIT301 > P3_P302
3550 11 ...
  
```

```

3551 10 =====
3552 11 P2_MV201 == 2.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3553   ↪ P3_LIT301 > min_P3_LIT301 + 24
3554 12 =====
3555 13 P2_P205 == P2_P203 == P2_MV201 == 2.0
3556 14 slope_P3_LIT301 == P2_P201
3557 15 P2_FIT201 > P2_MV201
3558 16 P2_FIT201 > P2_P201
3559 17 P2_FIT201 > P3_FIT301
3560 18 ...
3561 19 =====
3562 20 P3_P302 == 1.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3563   ↪ P3_LIT301 > min_P3_LIT301 + 24
3564 21 =====
3565 22 P2_P205 == P2_P203 == P2_MV201 == 2.0
3566 23 slope_P3_LIT301 == P3_P302 == P2_P201
3567 24 P2_FIT201 > P2_MV201
3568 25 P2_FIT201 > P3_FIT301
3569 26 ...
3570 27 =====
3571 28 P3_P302 == 2.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3572   ↪ P3_LIT301 > min_P3_LIT301 + 24
3573 29 =====
3574 30 P2_MV201 one of { 1.0, 2.0 }
3575 31 slope_P3_LIT301 one of { -1.0, 1.0 }
3576 32 P3_DPIT301 > P3_P302 > slope_P3_LIT301
3577 33 ...

```

Listing 6.9: Conditional Invariants for P2_MV201 and P3_P302

3578 Moreover, from the analysis of the invariants it becomes apparent that
 3579 when P3_P302 is in the ON state and P2_MV201 is in the OFF state, both
 3580 P3_FIT301 and P3_DPIT301 take values greater than 2 (as derived from
 3581 lines 5, 7, 8 and 32 of the listing).

3582 Regarding the valves P3_MV301 and P3_MV303, Daikon's analysis does
 3583 not provide specific information about their behavior during changes in
 3584 slope when the water level rises. Therefore, further analysis is required to
 3585 understand their role in the system.

3586 However, one observation can be made regarding P3_MV304. Daikon's
3587 analysis reveals two different slopes (increasing and decreasing) when this
3588 valve is in the ON state, which aligns with the observations made in the
3589 Graphical Analysis in Section 6.3.2.2. This finding strengthens the conjecture
3590 that P3_MV304 does not play a significant role in the tank cycle represented
3591 by P3_LIT301.

3592 **Analysis of the Current System Configuration** To simplify the analysis
3593 and facilitate the interpretation of outcomes, two separate groups of actuators
3594 were analyzed. The first group consists of P2_MV201 and P3_P302,
3595 which are conjectured to regulate the level of the tank. The second group
3596 includes P3_MV301, P3_MV303, and P3_MV304, for which it is speculated that
3597 they do not play a role in regulating the tank level. This grouping allows
3598 for a clearer examination of the states of the system and provides an opportunity
3599 to gain a better understanding of the behavior of these actuators.
3600 By analyzing these two groups separately, it becomes easier to draw conclusions
3601 and make comparisons between the different sets of actuators.

3602 The analysis of the first group of actuators, specifically P2_MV201 and
3603 P3_P302, further supports the conjectures made regarding their behavior
3604 in relation to the trend of the tank level. It was necessary to refine the
3605 analysis manually using the runDaikon.py script to obtain more detailed
3606 insights, particularly for the state [P2_MV201 == 2, P3_P302 == 2].

3607 Regarding the second group of actuators, the analysis of their states
3608 only reinforces the hypothesis that their activation does not have a significant
3609 impact on the trend of tank level represented by P3_LIT301.

3610 Unfortunately, no further useful information can be derived from this
3611 phase of the analysis. To address the remaining questions and clarify any
3612 outstanding issues, we will proceed to the next and final phase of the sub-
3613 system analysis.

3614 **6.3.2.4 Business Process Mining and Analysis**

3615 One of the hypotheses that needed to be tested was whether the valves
3616 P3_MV301, P3_MV303, and P3_MV304 have an impact on the tank level in the
3617 section between setpoints 850 and 900. Our initial assumption was that
3618 these actuators do not affect the level detected by P3_LIT301 because, as
3619 observed in the Graphical Analysis, the slope in that interval is similar
3620 to the slope between setpoints 970 and 1000, where these valves are not
3621 involved.

3622 To verify this conjecture, we utilized the JSON file generated by the
3623 processMining.py script. This script allows us to isolate the specific inter-
3624 vals and calculate the slope for each of them. In Listing 6.10, we present
3625 the outcomes of these calculations, providing further insights into the be-
3626 havior of the valves during those intervals.

```
3627 1 "slope_P3_LIT301": [  

3628 2 0.383, # from 970 to 1000  

3629 3 0.395, # from 850 to 900  

3630 4 0.384,  

3631 5 0.395,  

3632 6 0.354,  

3633 7 0.38,  

3634 8 0.388,  

3635 9 0.381,  

3636 10 0.385,  

3637 11 0.386  

3638 12 ] ,
```

Listing 6.10: Slope calculation of P3_LIT301 for the 850-900 and 970-1000 intervals related to tank levels

3639 The provided data in Listing 6.10 illustrates the calculated slopes for
3640 the intervals between 970 and 1000 (odd-numbered lines) and the inter-
3641 vals between 850 and 900 (even-numbered lines). Upon examination, we
3642 observe that the values for these two intervals are nearly identical, ac-
3643 counting for some expected fluctuations. This finding reinforces our initial
3644 conjecture that the P3_MV301, P3_MV303 and P3_MV304 valves do not play a

3645 role in the process of filling and emptying the tank. It is indeed possible to
3646 speculate that the P3_MV30x valves, similar to P3_MV302, might have a role
3647 in a different part of the system that has not been analyzed in the current
3648 context.

3649 Based on the activity diagram obtained from the analysis of actuators
3650 P2_MV201 and P3_P302 (Figure 6.9), the behavior of subsystem PLC2-3 can
3651 be summarized as follows:

- 3652 • When the system is in the states [P2_MV201 == 2, P3_P302 == 2] and
3653 [P2_MV201 == 2, P3_P302 == 1], the level of the tank represented by
3654 P3_LIT301 is increasing. The tank level exhibits a faster growth rate
3655 in the latter state.
- 3656 • When the system is in the state [P2_MV201 == 1, P3_P302 == 1], the
3657 tank level remains stable.
- 3658 • When the system is in the state [P2_MV201 == 1, P3_P302 == 2], the
3659 tank level is decreasing.

3660 The **absolute setpoints** for the tank level in subsystem PLC2-3 are de-
3661 fined as follows: the minimum setpoint is 800, the maximum setpoint is
3662 1000, and there are additional setpoints at 850, 900, and 970, which corre-
3663 spond to specific changes in the slope of the tank level.

3664 Taking a closer look at the behavior of sensors P2_FIT201 and P3_FIT301,
3665 we observe the following patterns: P2_FIT201, which is associated with
3666 incoming water flow and connected to valve P2_MV201, exhibits values
3667 greater than 2 when the valve is in the ON state. Conversely, when the
3668 valve is set to OFF, the sensor reading drops to 0.

3669 Regarding P3_FIT301, it appears to be linked to the behavior of the
3670 pump P3_P301 and correlates with the water flow out of the tank. When
3671 the pump is activated and in the ON state, the sensor records values greater
3672 than 2. On the other hand, when the pump is turned off and in the OFF
3673 state, the sensor reading returns to 0.

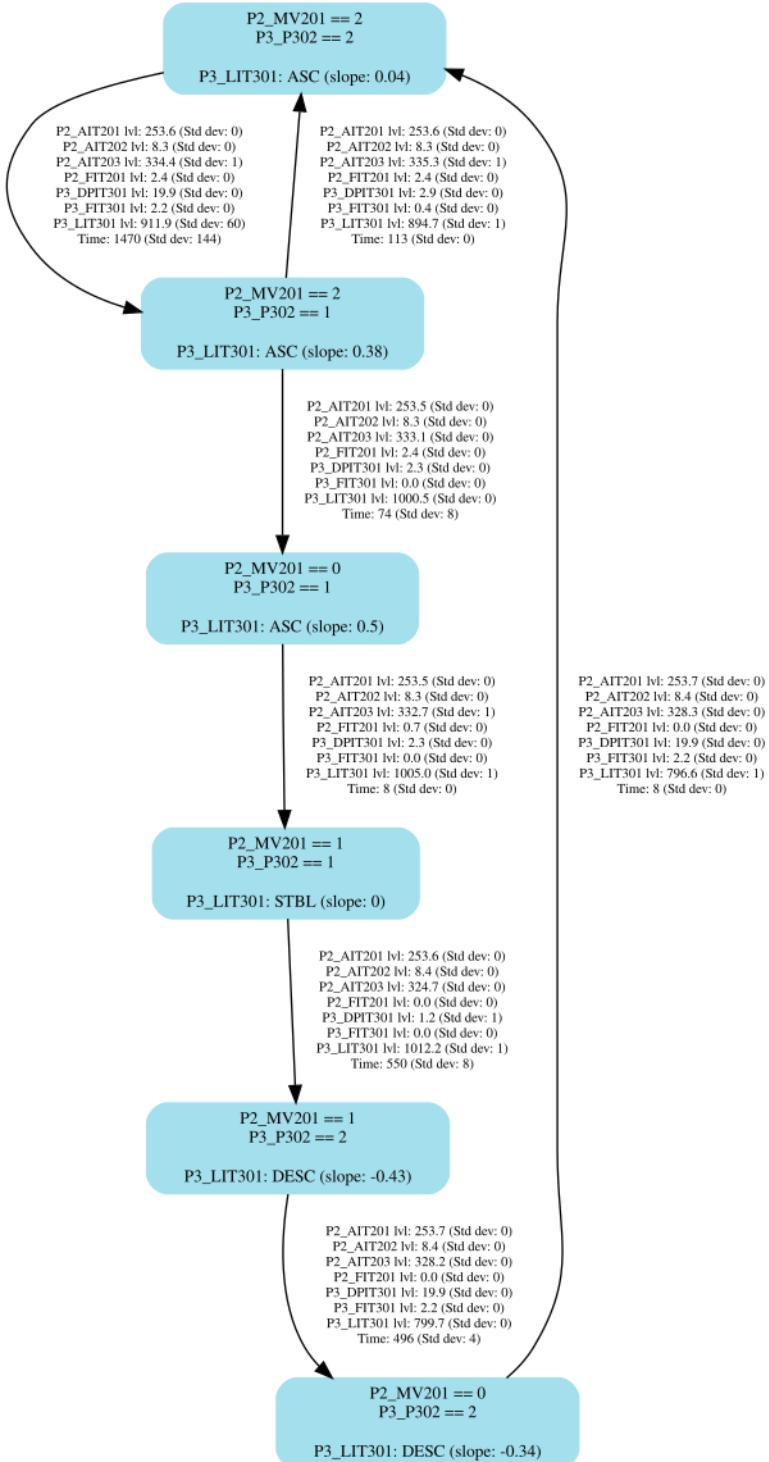


Figure 6.9: Activity diagram for PLC2-3

³⁶⁷⁴ **6.3.2.5 Properties**

³⁶⁷⁵ Table 6.2 provides a summary of the properties inferred from the con-
³⁶⁷⁶ jectures made throughout the different stages of the analysis.

#	Statement	Derived from
P23	The registers P3_DPIT301, P3_FIT301, and P3_LIT301 of PLC3 hold likely sensor measurements.	Preliminary Analysis Graphical Analysis
P24	The registers P3_MV301, P3_MV302, P3_MV303, P3_MV304, and P3_P302 of PLC3 hold likely actuator commands.	Preliminary Analysis Graphical Analysis
P25	The register P3_P301 of PLC3 is associated to a spare actuator.	Preliminary Analysis Graphical Analysis Invariant Analysis
P26	The register P3_LIT301 of PLC3 is associated to the level sensor of the tank controlled by PLC3.	Preliminary Analysis Graphical Analysis
P27	P2_MV201 and P3_P302 are associated to the actuators responsible for the level behavior of the water contained in the tank.	Graphical Analysis Invariant Analysis Business Process
P28	The rate of tank level growth is slow when both P2_MV201 and P3_P302 are in the ON state. The growth speed increases when P3_P302 switches to the OFF state.	Graphical Analysis Business Process
P29	When both P2_MV201 and P3_P302 are in the ON state, the inflow of water into the tank surpasses the outflow of water from the tank.	Graphical Analysis Business Process
P30	The tank level decreases when P2_MV201 is in the OFF state and P3_P302 is in the ON state.	Graphical Analysis Invariant Analysis
P31	The tank level remains (relatively) stable when both P2_MV201 and P3_P302 are in the OFF state.	Graphical Analysis Business Process

P32	The trend of the tank level switches from ascending to stable when the level reaches approximately 1000. It shifts from stable to descending when the level averages around 1012. It changes from descending back to ascending when the level reaches about 800. The slope of tank filling increases noticeably from around 850 to 900 and from around 970 to 1000.	Business Process
P33	Absolute setpoints are 800, 850, 900, 970, 1000.	Business Process
P34	P2_FIT201 is associated to a flow or pressure sensor to the P3_LIT301 register. It is related to the P2_MV201 valve.	Graphical Analysis Invariant Analysis Business Process
P35	P3_FIT301 and P3_DPIT301 exhibit similar patterns in their behavior. Both registers are related to the operation of pump P3_P302 and, consequently, to the flow of water out of the tank.	Graphical Analysis Business Process
P36	The registers P3_MV301, P3_MV303, and P3_MV304 do not have an impact on the water level dynamics of the tank controlled by PLC3.	Graphical Analysis Business Process

Table 6.2: Properties of the PLC2-3 subsystem

3677 6.3.3 Reverse Engineering of PLC3 and PLC4

3678 In the final phase of the reverse engineering process, the focus is di-
 3679 rected towards the subsystem consisting of PLC3 and PLC4 in the iTrust
 3680 SWaT system. Given the constraints of the thesis, this section will provide
 3681 a concise and schematic overview compared to the earlier sections.

3682 6.3.3.1 Pre-processing - Preliminary Analysis

3683 **Measurements and Actuators Recognition** Listing 6.11 shows the out-
 3684 comes obtained from automatic recognition of likely measurements and
 3685 actuators for PLC4. We omit those related to PLC3 as they are already
 3686 known.

3687 1 **Actuators:**
 3688 2 ...

```

3689   3
3690   4 Sensors:
3691   5 ...
3692   6 P4_AIT401 {'max_lvl': 148.8, 'min_lvl': 148.8}
3693   7 P4_AIT402 {'max_lvl': 191.1, 'min_lvl': 185.5}
3694   8 P4_FIT401 {'max_lvl': 1.7, 'min_lvl': 1.7}
3695   9 P4_LIT401 {'max_lvl': 1002.8, 'min_lvl': 775.8}
3696  10
3697  11 Hardcoded setpoints or spare actuators:
3698  12 ...
3699  13 P4_P401 [1.0]
3700  14 P4_P402 [2.0]
3701  15 P4_P403 [1.0]
3702  16 P4_P404 [1.0]
3703  17 P4_UV401 [2.0]

```

Listing 6.11: Preliminary analysis outcomes for sensors and actuators of PLC3-4

From the information provided in Listing 6.11, several observations can be made. Firstly, it is noted that there are no apparent actuators listed in the analysis. The likely sensors identified include P4_AIT401, P4_AIT402, P4_FIT401, and P4_LIT401. Among these sensors, P4_LIT401 is presumed to be the level sensor for the tank controlled by PLC4 based on similarities with the previous cases.

It is acknowledged that P4_FIT401 and P4_AIT401 are recognized as sensors, despite their seemingly constant values. It is important to note that the script used to identify likely actuators and sensors rounds the values to the first decimal place. Therefore, it is inferred that these registers contain continuous data with narrow value ranges.

Drawing on the analogy with the P21 property mentioned in Section 6.3.1.5, it is speculated that the P4_AIT40x registers represent measurements related to some water property. Additionally, P4_FIT401 is speculated to represent a pressure or flow sensor, based on similarities observed in previous cases.

In the analysis of the hardcoded setpoints and spare actuators, two registers stand out: P4_P402 and P4_UV401, both with a value of 2. Drawing on

3722 analogies from previous cases, P4_P402 is speculated to represent a pump
 3723 that is constantly in the ON state and therefore active. However, regarding
 3724 P4_UV401, it is unclear whether it is an actuator, a hardcoded setpoint, or a
 3725 different type of register. Further analysis is needed to determine its exact
 3726 purpose and functionality within the system.

3727 On the other hand, it can be concluded that P4_P401, P4_P403, and
 3728 P4_P404 are spare actuators, specifically pumps.

3729 **Actuator State Durations** Since there are no state-changing actuators within
 3730 PLC4, further analysis regarding the duration of actuator states will not
 3731 be performed for this subsystem. Please refer to Section 6.3.2.1 for evalua-
 3732 tions of the duration of actuator states in PLC3.

3733 **Actuator State Changes** As previously mentioned, our assumption is
 3734 that P4_LIT401 serves as the level sensor for the tank controlled by PLC4.
 3735 In our analysis of PLC2-3, we speculated that P3_MV302 acted as a valve
 3736 responsible for the incoming flow to an external element outside of that
 3737 subsystem. To test this hypothesis, we examine the setpoints of P3_MV302
 3738 in relation to the level of tank P4_LIT401. The setpoints of P3_MV302 corre-
 3739 sponding to the tank level are presented in Listing 6.12.

3740	1	Actuator state changes:		
3741	2	...		
3742	3	P3_MV302	prev_P3_MV302	P4_LIT401
3743	4	0	2	1000.5510
3744	5	0	1	784.4911
3745	6	0	2	922.1866
3746	7	0	1	881.0818
3747	8	0	2	1000.2820
3748	9	0	1	786.1061
3749	10	0	2	922.8018
3750	11	0	1	881.8508
3751	12

Listing 6.12: P3_MV302 state changes in relation to P4_LIT401

3752 The initial analysis suggests a potential correlation between the tank
3753 level values of P4_LIT401 and the behavior of P3_MV302. The ON state of
3754 P3_MV302 aligns with an increase in the tank level, while the OFF state cor-
3755 responds to a decrease. However, further analysis is required to provide
3756 additional evidence and support for this conjecture.

3757 **6.3.3.2 Graphs and Statistical Analysis**

3758 We will attempt to gain a deeper understanding of the pattern exhib-
 3759 ited by the PLC4 registers by referencing Figure 6.10.

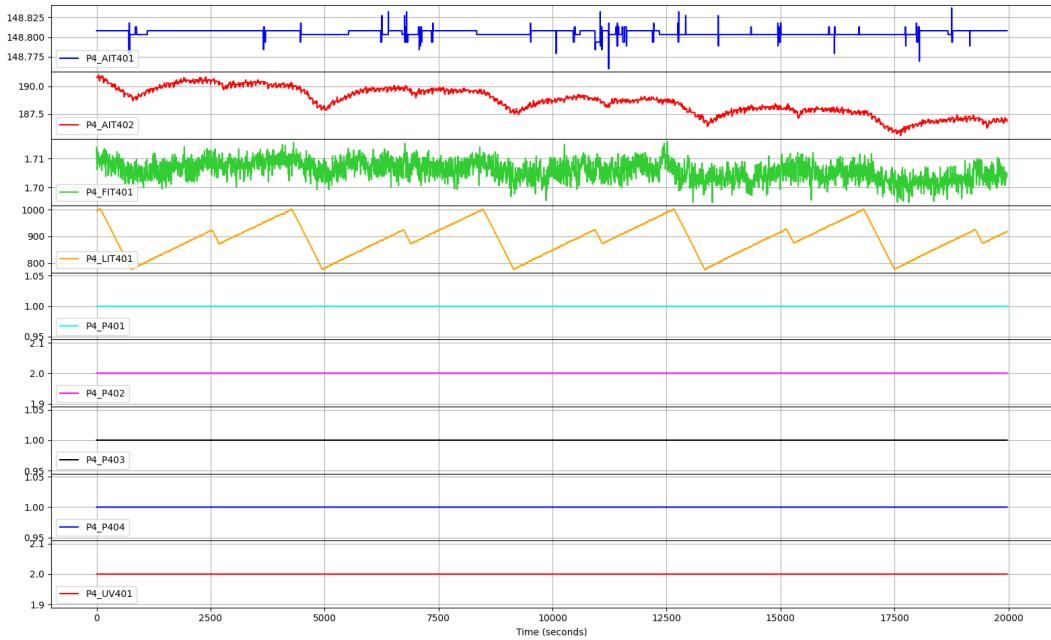


Figure 6.10: PLC4 registers

3760 The image reveals interesting behavior in the P4_AIT401 and P4_AIT401
 3761 registers. Notably, P4_AIT401 exhibits a linear trend rather than a cyclic
 3762 one, with values oscillating within a narrow range. This suggests that,
 3763 similar to P2_AIT201 (refer to Section 6.3.2.2), this register may correspond
 3764 to a sensor measuring a specific water property. On the other hand, P4_AIT402
 3765 appears to follow the level trend of sensor P4_LIT401, but with a down-
 3766 ward cyclic pattern where each cycle starts at a lower level than the pre-
 3767 vious one. Given the limited value range and the similarity in naming
 3768 conventions, it is highly likely that this register represents another water
 3769 property sensor rather than a tank.

3770 Additionally, P4_FIT401 does not display a cyclic pattern like the other
 3771 registers of the same type, and its values exhibit minimal variation, cor-
 3772 roborating the findings from the previous analysis phase.

3773 Now, let us refer to Figure 6.11 to seek confirmation regarding the con-
 3774 jecture that implicates valve P3_MV302 as the responsible actuator for tank
 3775 filling.

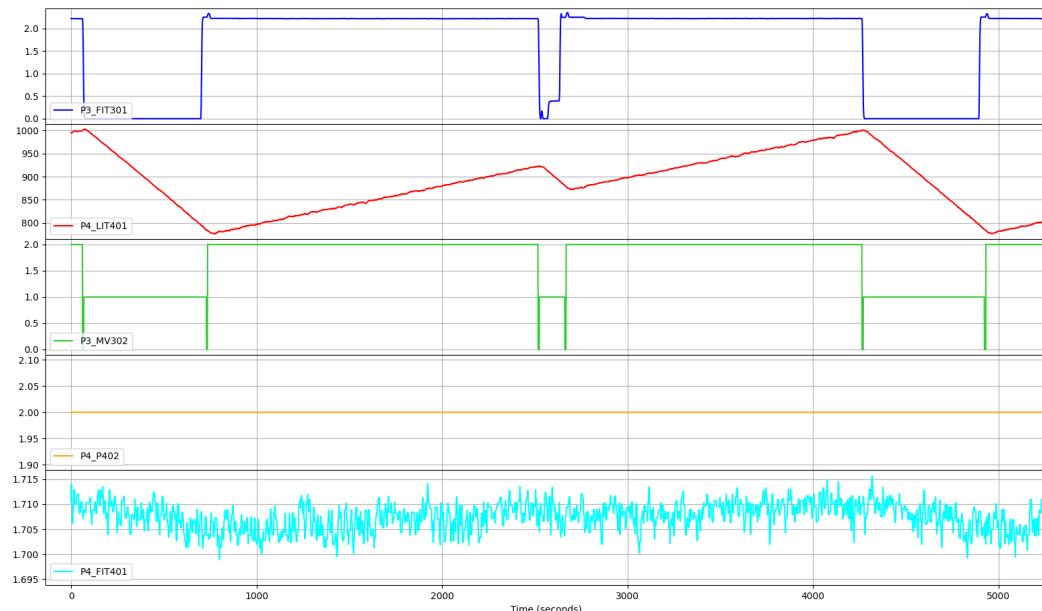


Figure 6.11: P3_MV302 and P4_LIT401 behaviors

3776 The image provides clear evidence that the behaviors of valve P3_MV302
 3777 and tank level sensor P4_LIT401 perfectly align. Additionally, P3_FIT301
 3778 (and its corresponding sensor, P3_DPIT301) appear to be related to the pat-
 3779 tern observed in P4_LIT401.

3780 Upon closer observation, it becomes apparent that the tank controlled
 3781 by PLC4 does **not have plateau periods**. When the incoming water flow
 3782 ceases, the tank immediately begins to empty. Based on our findings in
 3783 previous subsystems, we speculate that the actuator responsible for tank
 3784 emptying could be the pump indicated by register P4_P402. This specula-
 3785 tion is further supported by the nearly constant trend observed in sensor
 3786 P4_FIT401.

3787 Figure 6.12 depicts the correlation between the tanks within this sub-
 3788 system and the actuators that are responsible for their filling cycle.

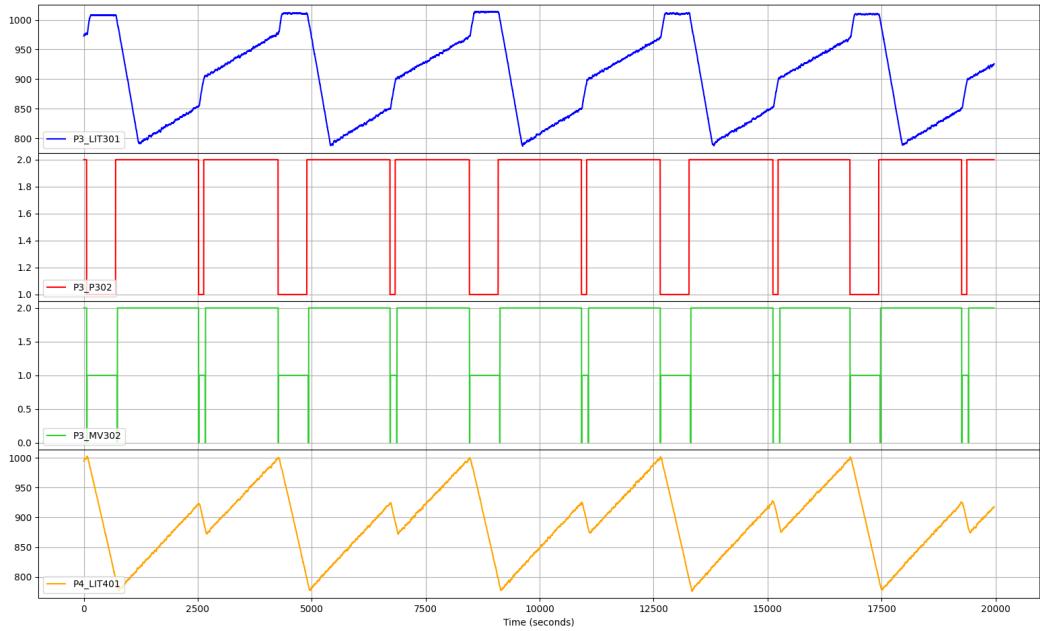


Figure 6.12: Tanks in subsystem PLC3-4 and their correlation.

3789 Further analysis was conducted to investigate whether valves P3_MV301,
 3790 P3_MV303, and P3_MV304 played a role in the tank filling cycle of PLC4.
 3791 However, the analysis did not confirm this hypothesis. Therefore, it can
 3792 be speculated that these valves are connected to other parts of the system
 3793 that are not currently discussed in this analysis.

3794 6.3.3.3 Invariant Inference and Analysis

3795 The invariant analysis for the subsystem consisting of PLC3-4 will be
 3796 brief as the states of the subsystem align with the states of valve P3_MV302,
 3797 with P4_P402 being constant throughout.

3798 **General Invariants** Again, the analysis of the general invariants offers
 3799 no new information compared to what we conjectured earlier. We there-
 3800 fore continue with the analysis of the current system configuration.

3801 **Analysis of the Current System Configuration** The analysis of the cur-
 3802 rent system configuration provides confirmation that when the system is

3803 in the state [$P3_MV302 == 1$, $P4_P402 == 2$], the tank level, as indicated
 3804 by sensor $P4_LIT401$, shows a decreasing slope. Additionally, it is noted
 3805 that the spare actuators align with the expected behavior.

3806 However, in the state [$P3_MV302 == 2$, $P4_P402 == 2$], the invariants
 3807 generated by Daikon do not provide the anticipated slope data, which was
 3808 expected to be increasing based on previous observations. This is due to
 3809 the descending slope between levels 880 and 925 of the tank represented
 3810 by $P4_LIT401$: by refining the analysis between the two increasing slope
 3811 sections we get the correct outcome, as shown in Listing 6.13.

```
3812 1 =====
3813 2 P3_MV302 == 2 && P4_P402 == 2 && P4_LIT401 < 990 &&
3814 3   ↪ P4_LIT401 > 930
3815 4 =====
3816 5 ...
3817 6 slope_P4_LIT401 == slope_P3_LIT301 == P4_P404 == P4_P403
3818 7   ↪ == P4_P401 == P3_P301 == P3_MV304 == P3_MV303 ==
3819 8   ↪ P3_MV301 == 1.0
3820 9 ...
3821 10 ...
```

Listing 6.13: Daikon manual analysis for $P3_MV302 == 2$

3821 6.3.3.4 Business Process Mining and Analysis

3822 The process mining step for the PLC3-4 subsystem is straightforward.
 3823 In this phase, we will not focus on determining the chronological order of
 3824 states (as we already know they correspond to the states of valve $P3_MV302$).
 3825 Instead, we will examine the setpoints and extract any available informa-
 3826 tion concerning additional measurements. Figure 6.13 presents the activity
 3827 diagram depicting the subsystem under analysis.

3828 The diagram provides confirmation that when the system is in the state
 3829 [$P3_MV302 == 2$, $P4_P402 == 2$], the tank level trend is increasing. Con-
 3830 versely, when the system is in the state [$P3_MV302 == 2$, $P4_P402 == 2$],
 3831 the trend is decreasing.

3832 Regarding the setpoints, based on the standard deviation of the data

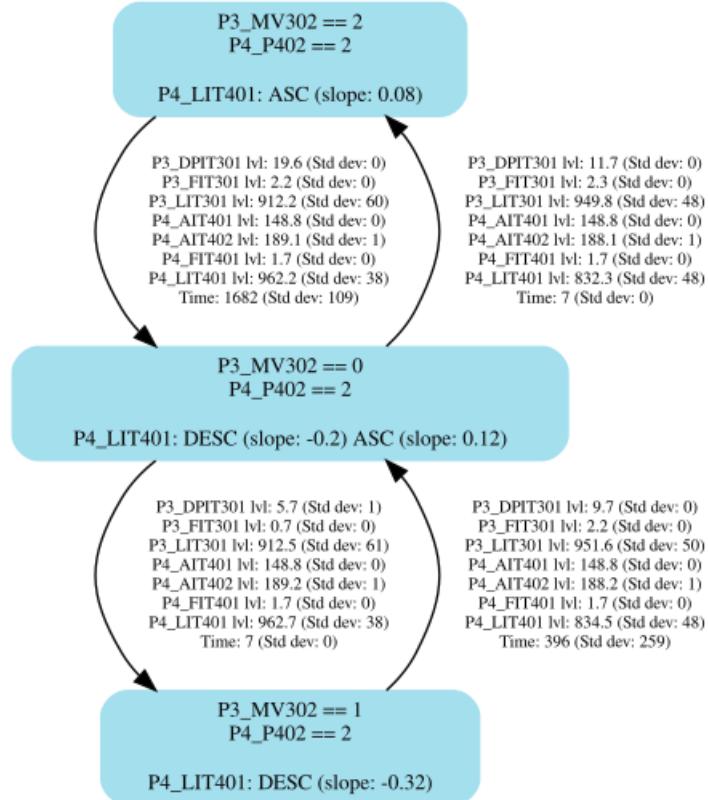


Figure 6.13: Activity diagram for PLC3-4

3833 from P4_LIT401, the absolute setpoints for the system are as follows: 785
 3834 (minimum), 880 (relative minimum), 925 (relative maximum), and 1000
 3835 (maximum).

3836 Furthermore, from the process mining phase, we can extract informa-
 3837 tion about P3_FIT301 and P3_DPIT301. P3_FIT301 registers values greater
 3838 than 2 when both pump P4_P402 and valve P3_MV302 are in the ON state,
 3839 while it drops to an average of 0.7 when the valve is off. Similarly, P3_DPIT301
 3840 registers values close to 20 when the valve is ON, and these values de-
 3841 crease when the valve is OFF.

3842 **6.3.3.5 Properties**

3843 Table 6.3 provides a summary of the properties inferred from the con-
3844 jectures made throughout the different stages of the analysis. Regrettably,
3845 we encountered difficulties in determining the type and purpose of the
3846 register labeled as P4_UV401.

#	Statement	Derived from
P37	The registers P4_AIT401, P4_AIT402, P4_FIT401, and P4_LIT401 of PLC4 hold likely sensor measurements.	Preliminary Analysis Graphical Analysis
P38	The register P4_P402 of PLC4 holds likely actuator command. Its status is constantly ON.	Preliminary Analysis Graphical Analysis Business Process
P39	The registers P4_P401, P4_P403, and P4_P404 of PLC4 are associated to spare actuators.	Preliminary Analysis Graphical Analysis Invariant Analysis
P40	The register P4_LIT401 of PLC4 is associated to the level sensor of the tank controlled by PLC4.	Preliminary Analysis Graphical Analysis
P41	P3_MV302 and P4_P402 are the actuators responsible for the level behaviour of the water contained in the tank controlled by PLC4.	Graphical Analysis Invariant Analysis Business Process
P42	The level of the tank identified by the register P4_LIT401 increases when both P3_MV302 and P4_P402 are in the ON state. It decreases when P3_MV302 is in the OFF state.	Graphical Analysis Invariant Analysis Business Process
P43	The trend of the tank level controlled by PLC4 transition from ascending to descending when the level reaches approximately 925 and 1000. It changes from descending when the level reaches approximately 785 and 880.	Business Process
P44	Absolute setpoints are 785, 880, 925 and 1000	Business Process
P45	P4_FIT401 serves as a flow or pressure sensor to the P4_LIT401 register. It is related to the P4_P402 pump.	Graphical Analysis Business Process
P46	P4_P401 does not exhibit a cyclic trend.	Graphical Analysis

P47	The register P4_P402 exhibits a cyclic decreasing trend, which is closely linked to the trend observed in the P4_LIT401 register.	Graphical Analysis
P48	Both P4_AIT401 and P4_AIT402 serve as sensors for measuring certain properties of the water.	Preliminary Analysis Graphical Analysis
P49	The registers P3_MV301, P3_MV303, and P3_MV304 do not have an impact on the water level dynamics on the tank controlled by PLC4.	Graphical Analysis Business Process

Table 6.3: Properties of the PLC3-4 subsystem

3847 6.4 PLCs Architecture

3848 In Table 6.4, we present an overview of the architecture of the four analyzed PLCs derived from the previous phases of the analysis. For each
 3849 PLC, we provide details about its constituent registers, including their op-
 3850 erational range. Where possible, we also indicate the role of these registers
 3851 within the physical system. By combining the information presented in
 3852 Tables 6.1, 6.2, and 6.3 with the content of Table 6.4, we can consider the
 3853 reverse engineering process implemented in this analysis to be complete.
 3854 To enhance readability, we will assign integer numbers to tanks based on
 3855 the PLCs that control them. For example, Tank 1 represents the tank asso-
 3856 ciated with PLC1, Tank 3 represents the tank associated with PLC3, and so
 3857 forth.

PLC	Physical device / Variable	Range	PLC register
PLC1	Level sensor for Tank 1	[489-815]	P1_LIT101
	Flow or pressure sensor	[0-2.7]	P1_FIT101
	Valve for Tank 1 inflow water	{0,1,2}	P1_MV101
	Pump for Tank 1 outflow water	{0,1}	P1_P101
	Spare actuator (OFF)	1	P1_P102
PLC2	Flow or pressure sensor	[0-2.4]	P2_FIT201
	Sensor for some water property	[252-256]	P2_AIT201
	Sensor for some water property	[8.3-8.4]	P2_AIT202

	Unkown	[320-342]	P2_AIT203
	Valve for Tank 3 inflow water	{0,1,2}	P2_MV201
	Spare actuator (OFF)	1	P2_P201
	Spare actuator (OFF)	1	P2_P202
	Pump (unknown role)	{0.1}	P2_P203
	Spare actuator (OFF)	1	P2_P204
	Pump (unknown role)	{0.1}	P2_P205
	Spare actuator (OFF)	1	P2_P206
PLC3	Level sensor for Tank 3	[768-1014]	P3_LIT301
	Flow or pressure sensor	[0-2.4]	P3_FIT301
	Flow or pressure sensor	[0-20.4]	P3_DPIT301
	Valve (unkonw role)	{0,1,2}	P3_MV301
	Valve for Tank 4 inflow water	{0,1,2}	P3_MV302
	Valve (unkonw role)	{0,1,2}	P3_MV303
	Valve (unkonw role)	{0,1,2}	P3_MV304
	Spare actuator (OFF)	1	P3_P301
	Pump for Tank 3 outflow water	[0,1]	P3_P302
	Level sensor for Tank 4	[775-1002]	P4_LIT401
PLC4	Flow or pressure sensor	[1.7]	P4_FIT401
	Sensor for some water property	[148.8]	P4_AIT401
	Sensor for some water property	[185-191]	P4_AIT402
	Spare actuator	1	P4_P401
	Pump for Tank 4 outflow water	2	P4_P402
	Spare actuator	1	P4_P403
	Spare actuator	1	P4_P404
	Unknown	2	P4_UV401

Table 6.4: Summary of the reconstructed composition PLCs from 1 to 4 in the iTrust SWaT System

Conclusion and Future Work

3859 **Discussion** At this stage, we have completed our (partial) reverse engineering
3860 of the iTrust SWaT System, aiming to achieve a sufficient level of
3861 the physical process, or *process comprehension*.

3862 Based on our analysis, we have identified a total of **49 properties** that
3863 are associated with the registers, their roles within the physical system,
3864 and their interactions with each other. These properties are summarized
3865 in Tables 6.1, 6.2, and 6.3. Additionally, we have reconstructed the **archi-**
3866 **tecture of the four examined PLCs**, providing details on the registers they
3867 consist of and their respective characteristics. This architecture, along with
3868 the specific registers associated with each PLC, is presented in Table 6.4.

3869 Nonetheless, despite the limited availability of network traffic data, we
3870 have managed to derive a basic schematic of the PLC network and the
3871 communication pathways between them. Please note that this schematic
3872 may be incomplete due to the lack of comprehensive network traffic data.
3873 The reader can refer to Figure 6.1 to visualize the network of PLCs.

3874 The obtained *process comprehension* enables us to plan a targeted attack
3875 on the system, utilizing the information we have gathered.

3876 To evaluate the accuracy of the information obtained about the SWaT
3877 system, we can refer to Figure 7.1 [70], which provides a schematic rep-
3878 resentation of the SWaT system. An \times indicates placement of sensors.

3879 By comparing the information derived from our analysis with the system
 3880 schematic, we can assess the validity and reliability of our findings.

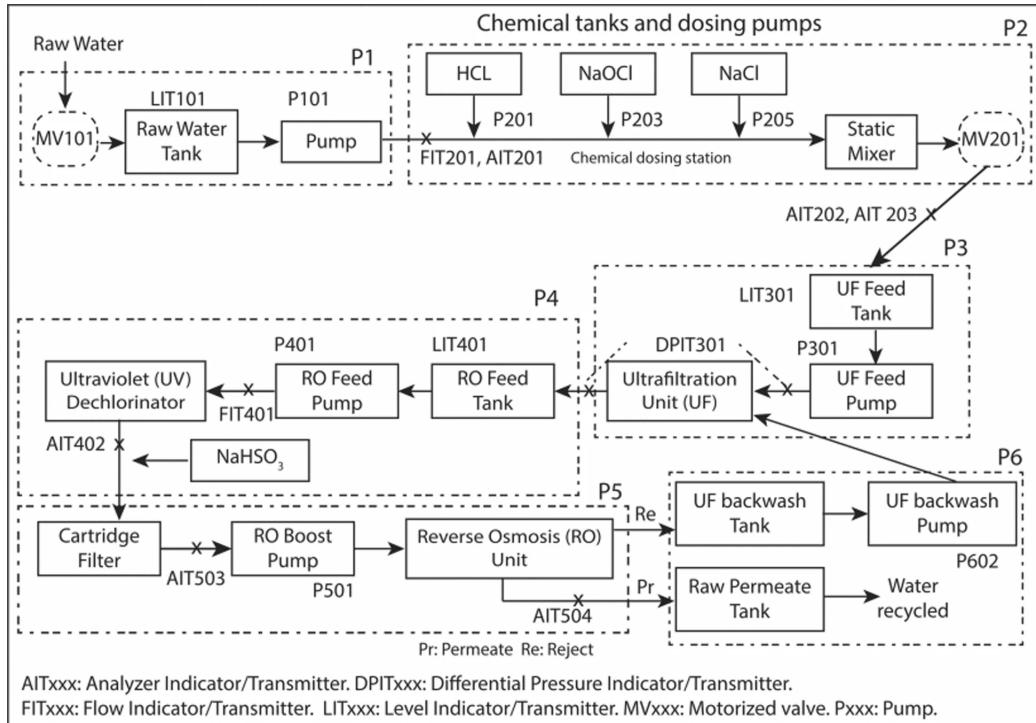


Figure 7.1: iTrust SWaT schema

3881 This image, while not comprehensive, serves to validate the accuracy
 3882 of the properties derived from our system analysis of the first four stages
 3883 of the SWaT system. It demonstrates that we have successfully identified
 3884 the actuators and sensors within the system, and in some cases, we have
 3885 even determined their specific roles within the physical process.

3886 Figure 7.2 [70] provides an alternative representation of the SWaT sys-
 3887 tem from the perspective of the Human-Machine Interface (HMI). This
 3888 depiction complements the previous diagram by adding more contextual
 3889 information and enhancing overall understanding of the system. It offers
 3890 a comprehensive view of the system's components and their relationships,
 3891 thereby improving the clarity and comprehension of the system's structure
 3892 and functioning.

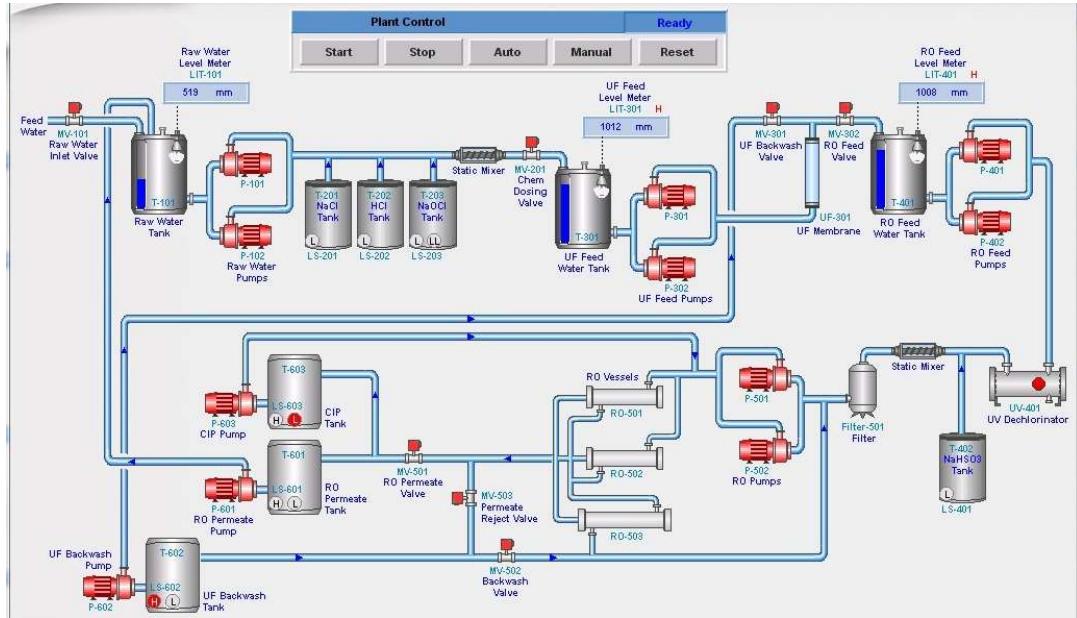


Figure 7.2: iTrust SWaT schema from HMI point of view.

3893 This figure introduces additional elements that were not included in
 3894 the previous schematic shown in Figure 7.1. It incorporates spare actuators
 3895 and other components such as chemical tanks and the ultrafiltration
 3896 membrane, which were not explicitly represented as registers in the avail-
 3897 able datasets. A comprehensive list of sensors and actuators, along with
 3898 their respective functions within the system, can be found in the paper by
 3899 Adepu et al. [73].

3900 **Achievements of the approach** After comparing the analysis results with
 3901 the actual system, we can assess the extent to which the proposed frame-
 3902 work has achieved its set objectives and identify any limitations encoun-
 3903 tered.

3904 In terms of the objectives that have been achieved, either fully or par-
 3905 tially, we have the following outcomes:

- 3906 • the proposed framework has effectively demonstrated independence
 3907 from the specific system being examined;

- 3908 • users have the ability to select a specific subsystem of their choice for
3909 analysis;
- 3910 • automatic recognition of registers associated with physical process
3911 measurements and actuators (even spare actuators) has been imple-
3912 mented, including their actual operating ranges (*relative setpoints*);
- 3913 • specific properties of actuators, including their states and the dura-
3914 tion of each state, have been identified;
- 3915 • recognition of which actuators handle incoming or outgoing flows;
- 3916 • partial recognition of the roles of certain registers within the system
3917 has been achieved, such as tanks, valves, pumps, and flow/pressure
3918 sensors;
- 3919 • conditions that cause state changes within the system (absolute set-
3920 points) have been recognized;
- 3921 • data perturbations for specific measurements have been attenuated;
- 3922 • the Graphical Analysis component now provides a comprehensive
3923 view of the system under investigation, enabling the extraction of
3924 much more information compared to the corresponding functional-
3925 ity of the Ceccato et al. framework;
- 3926 • the identification of invariants has been made easier, thanks in part
3927 to the implementation of semi-automated analyses;
- 3928 • the process mining capabilities related to the physical system can
3929 extract more data compared to the Ceccato et al. framework;
- 3930 • the reconstruction of the PLC network and the identification of the
3931 industrial protocol used for communications have been achieved.

3932 **Limitations of the approach** In terms of limitations, we have the follow-
3933 ing observations:

- 3934 • some specific register types or individual registers, such as valve
3935 P3_MV301, sensor P2_AIT203, or likely actuator P4_UV401, could not
3936 have their roles identified;
- 3937 • the presence and role of other system components, such as piping or
3938 filter membranes, could not be determined;
- 3939 • due to highly incomplete and partial available network traffic data,
3940 it was not possible to integrate network traffic data into the Business
3941 Process Analysis phase;
- 3942 • as a result of the aforementioned limitation, the Network Analysis
3943 phase is also incomplete, as it could not analyze communications
3944 occurring within the network via the CIP over EtherNet/IP protocol;
- 3945 • both the integration of network traffic data and the comprehensive
3946 analysis of network communications are considered as future work,
3947 to be addressed when more consistent and complete data become
3948 available.

3949 Based on our findings, we can confidently state that the proposed frame-
3950 work has successfully addressed many of the limitations identified in the
3951 Ceccato et al. framework. It has provided a greater amount of refined in-
3952 formation, aligning with the goals set forth at the beginning of this thesis.
3953 However, it is important to acknowledge that certain limitations, partic-
3954 ularly related to network data, have restricted the framework from fully
3955 realizing its potential. These limitations highlight the need for further ad-
3956 vancements in gathering comprehensive and complete network data to
3957 unlock the framework's full capabilities.

3958 **Future work** The framework, along with the thesis sources and analysis
3959 files, is openly accessible on the dedicated GitHub repository [75]. There is

3960 potential for further improvement and expansion of the framework. One
3961 possibility is integrating the analysis framework with an additional frame-
3962 work that utilizes the gathered data to generate targeted system attacks,
3963 building upon the writer's previous work as described in Section 3.2.2.

3964 Moreover, the proposed framework can benefit from various improve-
3965 ments and the introduction of specific new features to enhance its capa-
3966 bilities and effectiveness in analyzing industrial control systems. Here are
3967 some potential avenues for future work and enhancements to consider:

- 3968 • enhance the scanning and data gathering phase by reimplementing
3969 it to include support for multiple industrial protocols, in addition
3970 to Modbus. This improvement would enable the framework to de-
3971 tect and communicate with PLCs using various protocols commonly
3972 used in ICS environments, expanding its compatibility and usability;
- 3973 • enhance the automatic recognition of sensors and actuators by lever-
3974 aging advanced machine learning and artificial intelligence techniques.
3975 This improvement would involve training models to accurately iden-
3976 tify and classify different types of sensors and actuators based on
3977 their data patterns and characteristics;
- 3978 • complete the implementation of network traffic data within the Busi-
3979 ness Process part of the framework. However, it is crucial to ensure
3980 that the network traffic data and physical process data used for anal-
3981 ysis are reliable and not compromised by external attackers, as accu-
3982 rate and untampered data is essential for effective analysis;
- 3983 • implement an automated system for recognizing real system config-
3984 urations by excluding actuators that do not significantly contribute
3985 to changing the system behavior within the analyzed PLC. This au-
3986 tomatic recognition system can save time and effort by eliminating
3987 the need to manually filter out non-contributing actuators during the
3988 analysis process.

List of Figures

2.1	ICS architecture schema	12
2.2	PLC architecture	15
2.3	PLC communication schema	16
2.4	Comparison between ST language and Ladder Logic	17
2.5	Example of HMI for a water treatment plant	20
2.6	Modbus Request/Response schema	22
2.7	Modbus RTU frame and Modbus TCP frame	23
2.8	Example of packet sniffing on the Modbus protocol	25
2.9	OSI model for EtherNet/IP stack	26
3.1	Workflow of Ceccato et al.'s stages and operations with used tools	34
3.2	The simplified SWaT system used for running Ceccato et al. methodology	35
3.3	Output graphs from graph analysis	39
3.4	An example of Disco generated activity diagram for PLC2 .	43
3.5	Execution traces of InputRegisters_IW0 on the three PLCs .	45
3.6	Business process with states and Modbus commands for the three PLCs	47
3.7	Water physics compared: simulated physics in the Simulink model (a) and physics in a real system (iTrust SWaT) (b). Fluctuations in the tank level in (b), almost completely absent in (a), can be appreciated.	53

3.8	Behavior of the Graph Analysis on the Ceccato et al.'s tool	55
3.9	Histogram plot overshadowing statistical information shown on the terminal window in the background	56
3.10	Example of Daikon's output	58
4.1	Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination)	74
4.2	Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination) in textual mode	75
4.3	Slope comparison with granularity 30 (a), 60 (b) and 120 seconds (c)	83
4.4	Savitzky-Golay filter (blue line), STL decomposition (green) and RDP algorithm (orange) comparison	85
4.5	Slope after the application of the STL decomposition	86
4.6	The new slope representation (green line) and the smoothed measurement data obtaind with the STL decomposition (red)	87
4.7	Plotting registers on the same y-axis	91
4.8	Example of the new graph analysis	92
4.8	Example of the new graph analysis (cont.)	93
4.9	Example of transitive closure graphs for invariants in PLC1 of iTrust SWaT system	97
4.10	Directory (a) and outcome files (b) for the single actuator states analysis	101
4.11	Daikon outcome files for system configuration analysis. Each file represents a single system state	101
4.12	Activity diagram for PLC1 of the iTrust SWaT system	108
5.1	SWaT architecture	115
5.2	Sensors and actuators associated with each PLC	116
5.3	SWaT network architecture and zoning	118
5.4	SWaT stabilization	120
6.1	Simplified graph of the iTrust SWaT system network	125
6.2	Chart of PLC1-2 registers	132

6.3	Chart of PLC1-2 registers without spare actuators (particular)	133
6.4	Slope calculation anomaly (in the circle)	138
6.5	Activity diagram for PLC1-2	139
6.6	Verifying the conjecture about valve P2_MV201	148
6.7	PLC3 registers	149
6.8	P3_MV301 and P3_MV303 analysis	150
6.9	Activity diagram for PLC2-3	156
6.10	PLC4 registers	162
6.11	P3_MV302 and P4_LIT401 behaviors	163
6.12	Tanks in subsystem PLC3-4 and their correlation.	164
6.13	Activity diagram for PLC3-4	166
7.1	iTrust SWaT schema	172
7.2	iTrust SWaT schema from HMI point of view.	173

List of Tables

2.1	Differences between Information Technology (IT) and Industrial Control Systems (ICSs)	10
2.2	Modbus Function Codes list	24
3.1	Summary table of Ceccato et al. framework limitations	60
4.1	Summary table of proposed framework enhancements	112
5.1	Main IP addresses of the six PLCs and SCADA in the SWaT system	118
5.2	SWaT network traffic data	121
6.1	Properties of the PLC1-2 subsystem	143
6.2	Properties of the PLC2-3 subsystem	158
6.3	Properties of the PLC3-4 subsystem	168
6.4	Summary of the reconstructed composition PLCs from 1 to 4 in the iTrust SWaT System	169

Listings

3.1	Example of registers capture	37
3.2	Example of raw network capture	38
3.3	Statistical data for PLC1_InputRegisters_IW0 register	40
3.4	Generic example of a .spinfo file for customizing rules in Daikon	41
3.5	The three sections of Daikon analysis outcomes	41
4.1	Novel Framework structure and Python scripts	65
4.2	Paths and parameters for the Pre-processing phase in <i>config.ini</i> file	77
4.3	config.ini parameters for dataset enriching	79
4.4	Example of preliminary system analysis	88
4.5	Standard Daikon output for PLC1 of the iTrust SWaT system	95
4.6	Revised Daikon output with transitive closures for PLC1 of the iTrust SWaT system	98
4.7	Daikon outcomes for the system configuration MV101 == 2, P101 == 1 on LIT101	102
4.8	Example of the data contained in the produced JSON file	106
6.1	Preliminary analysis outcomes for sensors and actuators of PLC1-2	128
6.2	Time duration of the states of actuators P1_MV101 and P1_MV201 of PLC1-2	130
6.3	P1_P101 state changes in relation to P1_LIT101	130
6.4	General Invariants for PLC1-2	135

6.5	Conditional Invariants for states 1 and 2 of P1_MV101	136
6.6	Preliminary analysis outcomes for sensors and actuators of PLC2-3	143
6.7	Time duration of the states of actuators of PLC3	144
6.8	P2_MV201 state changes in relation to P3_LIT301	146
6.9	Conditional Invariants for P2_MV201 and P3_P302	151
6.10	Slope calculation of P3_LIT301 for the 850-900 and 970-1000 intervals related to tank levels	154
6.11	Preliminary analysis outcomes for sensors and actuators of PLC3-4	158
6.12	P3_MV302 state changes in relation to P4_LIT401	160
6.13	Daikon manual analysis for P3_MV302 == 2	165

References

- [1] R. Rajkumar, L. Lee, I. Sha, and J. A. Stankovic. "Cyber-physical systems:the next computing revolution." In: 47th Design Automation Conference, DAC 2010, Anaheim, California, USA (July 2010). 2010. DOI: <http://dx.doi.org/10.1145/1837274.1837461>.
- [2] NIST. *ICS defintion*. URL: https://csrc.nist.gov/glossary/term/industrial_control_system (visited on 2023-04-06).
- [3] A. Akbarzadeh. "Dependency based risk analysis in Cyber-Physical Systems". PhD thesis. Norwegian University of Science, Technology - Faculty of Information Technology, and Electrical Engineering, 2023. ISBN: 978-82-326-6986-8 (electronic version).
- [4] Schneider Electrics. *What is Modbus and how does it work?* URL: <https://www.se.com/us/en/faqs/FA168406/> (visited on 2023-04-13).
- [5] Wikipedia. *DNP3*. URL: <https://en.wikipedia.org/wiki/DNP3> (visited on 2023-05-18).
- [6] *Introduction to EtherNet/IP Technology*. URL: https://scadahacker.com/library/Documents/ICS_Protocols/Intro%20to%20EthernetIP%20Technology.pdf (visited on 2023-04-19).
- [7] OPC Foundation. *Unified Architecture*. URL: <https://opcfoundation.org/about/opc-technologies/opc-ua/> (visited on 2023-06-15).
- [8] Wireshark.org. *S7comm*. URL: <https://wiki.wireshark.org/S7comm> (visited on 2023-06-15).

- [9] B. Green, M. Krotofil, and A. Abbasi. “On the Significance of Process Comprehension for Conducting Targeted ICS Attacks”. In: Workshop on Cyber-Physical Systems Security and PrivaCy (Nov. 2017). ACM, 2017, pp. 57–67. DOI: <https://doi.org/10.1145/3140241.3140254>.
- [10] ODVA Inc. “EtherNet/IP - CIP on Ethernet Technology”. In: URL: https://www.odva.org/wp-content/uploads/2021/05/PUB00138R7_Tech-Series-EtherNetIP.pdf (visited on 2023-04-18).
- [11] K. Pal, A. Adepu, and J. Goh. “Cyber-Physical System Discovery: Reverse Engineering Physical Processes”. In: 3rd ACM Workshop on Cyber-Physical System Security (Apr. 2017). 2017, pp. 3–14. DOI: <https://doi.org/10.1145/3055186.3055195>.
- [12] M. Ceccato, Y. Driouich, R. Lanotte, M. Lucchese, and M. Merro. “Towards Reverse Engineering of Industrial Physical Processes”. In: CPS4CIP@ESORICSAt 2022 (Copenhagen, Denmark, Sept. 30, 2022). 2022, 273–290. DOI: https://doi.org/10.1007/978-3-031-25460-4_15.
- [13] A. Keliris and M. Maniatakos. “ICSREF: A Framework for Automated Reverse Engineering of Industrial Control Systems Binaries”. In: 26th Annual Network and Distributed System Security Symposium (NDSS) (Feb. 2019). 2019. DOI: <https://doi.org/10.48550/arXiv.1812.03478>.
- [14] *The Daikon dynamic invariant detector*. URL: <http://plse.cs.washington.edu/daikon/> (visited on 2023-04-25).
- [15] iTrust Singapore University of Technology and Design. *Secure Water Treatment*. URL: <https://itrust.sutd.edu.sg/testbeds/secure-water-treatment-swat/> (visited on 2023-05-15).
- [16] Wikipedia. *Stuxnet*. URL: <https://en.wikipedia.org/wiki/Stuxnet> (visited on 2023-05-18).
- [17] Wikipedia. *Industroyer*. URL: <https://en.wikipedia.org/wiki/Industroyer> (visited on 2023-05-18).

- [18] F. Robles and N. Perlroth. “‘Dangerous Stuff’: Hackers Tried to Poison Water Supply of Florida Town”. In: *The New York Times* (2021-02-08). URL: <https://www.nytimes.com/2021/02/08/us/oldsmar-florida-water-supply-hack.html> (visited on 2023-05-18).
- [19] G. M. Makrakis, C. Koliас, G. Kambourakis, C. Rieger, and J. Benjamin. “Industrial and Critical Infrastructure Security: Technical Analysis of Real-Life Security Incidents”. In: *IEEE Access* (2021-12-06). DOI: <https://dx.doi.org/10.1109/ACCESS.2021.3133348>.
- [20] NIST. *PLC defintion*. URL: https://csrc.nist.gov/glossary/term/programmable_logic_controller (visited on 2023-04-06).
- [21] W. Bolton. *Programmable Logic Controllers, 6th edition*. Newnes, 2015, pp. 7–9.
- [22] H. S. G. Pussewalage, P. S. Ranaweera, and V. Oleshchuk. “PLC security and critical infrastructure protection”. In: 2013 IEEE 8th International Conference on Industrial and Information Systems (Dec. 2013). 2013. DOI: <https://dx.doi.org/10.1109/ICIIInfS.2013.6731959>.
- [23] Wikipedia. *Remote Terminal Unit*. URL: https://en.wikipedia.org/wiki/Remote_terminal_unit (visited on 2023-04-08).
- [24] *What is SCADA?* URL: <https://www.inductiveautomation.com/resources/article/what-is-scada> (visited on 2023-04-05).
- [25] NIST. *HMI defintion*. URL: https://csrc.nist.gov/glossary/term/human_machine_interface (visited on 2023-04-11).
- [26] Zscaler. *What Is the Purdue Model for ICS Security?* URL: <https://www.zscaler.com/resources/security-terms-glossary/what-is-purdue-model-ics-security> (visited on 2023-05-17).
- [27] O. Morando. *Hacking: Modbus - Cyber security OT/ICS*. URL: <https://blog.omarmorando.com/hacking/ot-ics-hacking/hacking-modbus> (visited on 2023-04-15).

- [28] Modbus.org. "MODBUS/TCP Security". In: pp. 7–8. URL: https://modbus.org/docs/MB-TCP-Security-v21_2018-07-24.pdf (visited on 2023-04-16).
- [29] Odva, Inc. URL: <https://www.odva.org> (visited on 2023-04-21).
- [30] ODVA Inc. *Common Industrial Protocol*. URL: <https://www.odva.org/technology-standards/key-technologies/common-industrial-protocol-cip/> (visited on 2022-12-26).
- [31] Wikipedia. *Common Industrial Protocol*. URL: https://en.wikipedia.org/wiki/Common_Industrial_Protocol (visited on 2023-04-21).
- [32] *What is the Common Industrial Protocol (CIP)?* URL: <https://www.motioncontrolltips.com/what-is-the-common-industrial-protocol-cip/> (visited on 2023-04-21).
- [33] *CIP (Common Industrial Protocol): CIP messages, device types, implementation and security in CIP*. URL: <https://resources.infosecinstitute.com/topic/cip/> (visited on 2023-04-21).
- [34] Wikipedia. *CODESYS*. URL: <https://en.wikipedia.org/wiki/CODESYS> (visited on 2023-05-18).
- [35] Y. Yuan, X. Tang, W. Zhou, W. Pan, X. Li, H. T. Zhang, H. Ding, and J. Goncalves. "Data driven discovery of cyber physical systems". In: *Nature Communications* 10 (2019-10-25). URL: <https://www.nature.com/articles/s41467-019-12490-1> (visited on 2023-04-20).
- [36] C. Feng, V.R. Palleti, A. Mathur, and D. Chana. "A Systematic Framework to Generate Invariants for Anomaly Detection in Industrial Control Systems". In: NDSS Symposium 2019 (San Diego, CA, USA, Feb. 24–27, 2019). 2019. DOI: <https://dx.doi.org/10.14722/ndss.2019.23265>.
- [37] Singapore University of Technology and Design. *iTrust - Center for Research in Cyber Security*. URL: <https://itrust.sutd.edu.sg/> (visited on 2023-05-15).

- [38] K. Pal, A. Adepu, and J. Goh. "Effectiveness of Association Rules Mining for Invariants Generation in Cyber-Physical Systems". In: 2017 IEEE 18th International Symposium on High Assurance Systems Engineering (HASE) (Jan. 2017). 2017. DOI: <https://doi.org/10.1109/HASE.2017.21>.
- [39] *Association rules mining*. URL: https://en.wikipedia.org/wiki/Association_rule_learning (visited on 2023-04-21).
- [40] B. Perelman. "Reconnaissance in Industrial Networks: What You Don't See Can Hurt You". In: *SecurityWeek* (2017-04-04). URL: <https://www.securityweek.com/reconnaissance-industrial-networks-what-you-dont-see-can-hurt-you/> (visited on 2023-05-18).
- [41] *Ceccato et al. reverse engineering prototype tool*. URL: <https://github.com/SeUniVr/PLC-RE> (visited on 2023-04-23).
- [42] R. Lanotte, M. Merro, and A. Munteanu. "Industrial Control Systems Security via Runtime Enforcement". In: *ACM Transactions on Privacy and Security* 26.4 (2023-02), pp. 1–41. DOI: <https://doi.org/10.1145/3546579>.
- [43] M. Oliani. *AttackPLC - Project for Network Security teaching*. URL: <https://github.com/marcooliani/AttackPLC> (visited on 2023-04-23).
- [44] *Ray - Productionizing and scaling Python ML workloads simply*. URL: <https://www.ray.io/> (visited on 2023-04-25).
- [45] *Tshark*. URL: <https://tshark.dev/> (visited on 2023-04-25).
- [46] *Wireshark*. URL: <https://www.wireshark.org/> (visited on 2023-04-25).
- [47] *pandas - Python Data Analysis library*. URL: <https://pandas.pydata.org/> (visited on 2023-04-25).
- [48] *NumPy - The fundamental package for scientific computing with Python*. URL: <https://numpy.org/> (visited on 2023-04-25).
- [49] *R project for statistical computing*. URL: <https://www.r-project.org/> (visited on 2023-04-25).

- [50] *SciPy - Fundamental alghorithms for scientific computing with Python.* URL: <https://scipy.org/> (visited on 2023-04-25).
- [51] *Enhancing Daikon output.* URL: <https://plse.cs.washington.edu/daike/download/doc/daike/Enhancing-Daikon-output.html#Splitter-info-file-format> (visited on 2023-04-25).
- [52] *Process mining.* URL: https://en.wikipedia.org/wiki/Process_mining (visited on 2023-04-26).
- [53] *Fluxicon Disco.* URL: <https://fluxicon.com/disco/> (visited on 2023-04-26).
- [54] *OpenPLC - Open source PLC software.* URL: <https://openplcproject.com/> (visited on 2023-04-23).
- [55] *Docker.* URL: <https://www.docker.com/> (visited on 2023-04-26).
- [56] MathWorks. *Simulink.* URL: <https://it.mathworks.com/products/simulink.html> (visited on 2023-04-26).
- [57] T. Carlsson. "Industrial network market shares 2023". In: *HMS Networks* (2023-05-05). URL: <https://www.hms-networks.com/news-and-insights/news-from-hms/2023/05/05/industrial-network-market-shares-2023> (visited on 2023-05-18).
- [58] *ProM Tools - The Process Mining Framework.* URL: <https://promtools.org/> (visited on 2023-04-26).
- [59] Fraunhofer Institute for Applied Information Technology. *pm4py - Process Mining for Python.* URL: <https://promtools.org/> (visited on 2023-04-26).
- [60] statsmodels. *statsmodels - statistical models, hypothesis tests, and data exploration.* URL: <https://www.statsmodels.org/stable/index.html> (visited on 2023-05-05).
- [61] NetworkX. *NetworkX - Network Analysis in Python.* URL: <https://networkx.org/> (visited on 2023-05-05).
- [62] Wikipedia. *Polynomial regression.* URL: https://en.wikipedia.org/wiki/Polynomial_regression (visited on 2023-05-21).

- [63] Wikipedia. *Decomposition of time series*. URL: https://en.wikipedia.org/wiki/Decomposition_of_time_series (visited on 2023-05-21).
- [64] M. Fleischmann. "Line simplification algorithms". In: (2020-04-27). URL: <https://martinfleischmann.net/line-simplification-algorithms/> (visited on 2023-05-21).
- [65] Wikipedia. *Savitzky–Golay filter*. URL: https://en.wikipedia.org/wiki/Savitzky-Golay_filter (visited on 2023-05-06).
- [66] Otext. *STL decomposition*. URL: <https://otexts.com/fpp2/stl.html> (visited on 2023-05-06).
- [67] Wikipedia. *Ramer–Douglas–Peucker algorithm*. URL: https://en.wikipedia.org/wiki/Ramer-Douglas-Peucker_algorithm (visited on 2023-05-21).
- [68] *The Daikon Invariant Detector User Manual*. URL: <https://plse.cs.washington.edu/daikon/download/doc/daikon.html#Program-point-sections> (visited on 2023-05-25).
- [69] Wikipedia. *Transitive closures*. URL: https://en.wikipedia.org/wiki/Transitive_closure (visited on 2023-05-11).
- [70] iTrust Singapore University of Technology and Design. *Secure Water Treatment (SWaT) Testbed*. URL: https://itrust.sutd.edu.sg/wp-content/uploads/2021/07/SWaT_technical_details-160720-v4.4.pdf (visited on 2023-05-15).
- [71] A. Mathur and N. O. Tippenhauer. "SWaT: a water treatment testbed for research and training on ICS security". In: 2016 International Workshop on Cyber-physical Systems for Smart Water Networks (CySWater) (Apr. 2016). 2016. DOI: <http://dx.doi.org/10.1109/CySWater.2016.7469060>.
- [72] iTrust Singapore University of Technology and Design. *Datasets*. URL: https://itrust.sutd.edu.sg/itrust-labs_datasets/ (visited on 2023-05-15).

- [73] S. Adepu, J. Goh, K.N. Junejo, and A. Mathur. "A Dataset to Support Research in the Design of Secure Water Treatment Systems". In: The 11th International Conference on Critical Information Infrastructures Security (France, Oct. 2016). 2016.
- [74] T. A. Snide. *CIP-Modbus Integration*. 2008. URL: https://www.modbus.org/docs/CIP%20Modbus%20Integration%20Hanover%20Fair_0408.pdf (visited on 2023-05-17).
- [75] M. Oliani. *Master Degree's Thesis Framework and Sources*. URL: <https://github.com/marcooliani/Thesis> (visited on 2023-06-15).